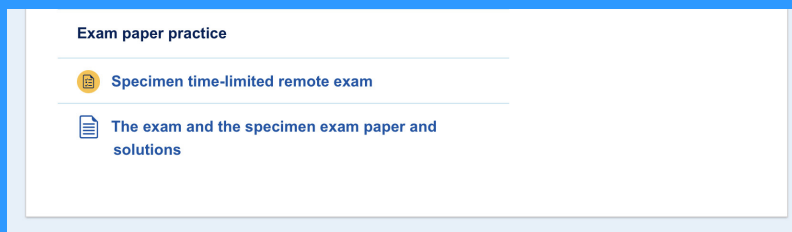


Welcome to this run through 2020.
download from: apina.dev/M249

There's a specimen remote exam at the assessment page of the course website, which is great to do a 'dry run' and practice your setup, scanning and uploading handwritten answers etc.



Aim to spend about 45 minutes on each section, so if a question looks gnarly and time consuming, skip it for now, stay positive and come back later



M2492006F1PV1



M249/J

Module Examination 2020
Practical Modern Statistics

Friday 12 June 2020

10.00 am – 1.00 pm

Time allowed: 3 hours

There are **four sections** in this examination, each worth 25% of the total mark.

In **each section** you should **attempt all questions**. You are advised to spend about 45 minutes on each section.

Include all your working, as some marks are awarded for this.

Write your answers in the spaces provided in this question paper in **pen**, though you may draw diagrams in pencil.

Crossed out work will not be marked.

You may request an answer book for further work, if needed.

Please fill in the grid below to show your Personal Identifier and Examination number, taken from your desk record.

Personal Identifier								
Examination No.								
Calculator Model								

For Examiner's use only:

Section	1	2	3	4	Total
Marks					

Section 1 (relates to Book 1 *Medical statistics*)

Questions 1 to 6

You should **attempt all questions**. *This section is worth 25%*.

Write your answers in the spaces provided.

Question 1

A study was undertaken to investigate whether an elevated level of an inflammatory marker is associated with an increased risk of coronary heart disease. Blood samples were taken from a sample of 18 225 men and their levels of the inflammatory marker were determined. Over the following three years, 266 of the men had a heart attack and their levels of the inflammatory markers were compared with those of a subset of the other men. That subset was chosen to match the men with a heart attack in terms of age and smoking status, with two men without a heart attack selected for each man who had suffered a heart attack.

- (a) State whether this is a cohort or case-control study, and give a reason for your answer.

[2]

Case-control

Group of cases with disease (infarction) is then matched with comparable cases without disease (investigates chose 2x number).

[To match in a cohort study, we would do this before the exposure, e.g. matching high marker with low marker and seeing how many in each group then had a heart attack]

- (b) In the study, each man who had suffered a heart attack was matched for age and smoking status with two men who had not. What bias could this remove? Explain briefly how this bias might have arisen without the matching.

[2]

Confounder bias.

Eg. if marker are higher for older people, as age is also associated with increased heart attack risk then age can be a confounder.

Question 2

A potential side effect of a medicine is an upset stomach. In a trial, 56 volunteers received a new medicine and a further 85 volunteers received a placebo. The numbers having an upset stomach are shown in Table 1.

Table 1

	Upset Stomach		Total
	Yes	No	
New medicine	14	42	56
Placebo	2	83	85
Total	16	125	141

- (a) Calculate the relative risk of an upset stomach for people taking the new medicine, relative to those taking the placebo.

[1]

HS PAGE 10, POINT 5

$$\frac{14/56}{2/85} = 10.625$$

- (b) Obtain a 95% confidence interval for the relative risk.

[2]

z quantile is 1.96.

$$\hat{\sigma} = \sqrt{\frac{1}{14} - \frac{1}{56} + \frac{1}{2} - \frac{1}{85}} = 0.7361$$

$$\text{so } 95\% \text{ CI is } (10.625 \exp\{-z\sigma\}, 10.625 \exp\{z\sigma\}) \\ = (2.51, 44.97)$$

(c) Summarize and interpret your results.

[2]

RR is 10.625, 95% CI (2.51, 44.97).

As confidence interval does not contain 1, the result is significant at the 95% confidence interval.

There is evidence that upset stomachs are more common with the new medicine than with placebo.

Question 3

A study of association between obesity and cardiovascular disease (CVD) was undertaken. The data were split into two strata according to age (Stratum A: age < 50 years; Stratum B: age \geq 50 years). The stratified data are shown in Table 2.

Table 2

A: Age under 50

Exposure category	CVD	No CVD
Obese	10	90
Not obese	35	465

B: Age 50 or more

Exposure category	CVD	No CVD
Obese	36	164
Not obese	25	175

- (a) Obtain the stratum specific odds ratios for the two groups and the odds ratio for the aggregated data.

[2]

HB PAGE 11, POINT 6

$$OR = \frac{ad}{bc}$$

$$\text{so } OR_{<50} = 1.4762$$

$$\& OR_{\geq 50} = 1.5366$$

$$\& OR_{\text{COMBINED}} = 1.9318$$

(b) Calculate the Mantel-Haenszel odds ratio

[2]

$$OR_{MH} = \frac{\sum a_i d_i / N_i}{\sum b_i c_i / N_i} \quad N_{i1} = a_i + b_i + c_i + d_i$$
$$N_1 = 600 \quad N_2 = 400$$

$$OR_{MH} = \frac{(10 \times 465) / 600 + (36 \times 175) / 400}{(90 \times 35) / 600 + (164 \times 25) / 400} = 1.5161$$

(c) Comment on the association between obesity and CVD and the role played by age in this association.

[2]

The OR for each strata 1.476 age < 50 & 1.537 together with the combined ORs are all > 1, suggesting obesity is positively associated with CVD.

The Mantel 1.5161 is between the stratum-specific OR; unlike the aggregated OR 1.9318 which suggests age is a confounder, with the smaller age > 50 category having a considerably greater proportion in the "obese" exposure category, age seems associated with exposure.

Question 4

A case-control study was conducted to examine whether there is an association between dietary magnesium intake and risk of colorectal tumors. The intake level was divided into three ranges that each contained similar numbers of controls. The dose-specific odds ratios relative to individuals in the lowest range are shown in Table 3.

Table 3

Magnesium intake	Cases	Controls	Odds ratio
Under 316	270	236	1.00
316 to 358	267	236	0.90
Over 358	231	237	0.73

- (a) Briefly describe the relationship between magnesium intake and colorectal tumors. [1]

Increasing magnesium intake is associated with decreasing odds ratio and hence lower risk of colorectal tumor.

- (b) The chi-squared test for no linear trend gives the value 1.523 for the test statistic. What do you conclude, and why? [2]

The null distribution for the test statistic is $\sim \chi^2(1)$.
(HB PAGE 13, POINT 18)

Quantiles for $\chi^2(1)$ are given on the first row of table 3 HB page 27

Note that $z_{0.8} = 1.64$ so here $p > 0.2$

Hence there is little evidence to reject the null-hypothesis of no linear trend ($\beta=0$) and we cannot infer a linear trend between increasing magnesium intake & colorectal cancer.

Question 5

A randomized controlled trial is planned to investigate whether a new drug for treating arthritis is better than a standard treatment. The investigators are interested in the proportion of arthritis sufferers who have reduced symptoms after three weeks. The sample size for the trial was based on the following values at significance level $\alpha = 0.01$, power $\gamma = 0.85$ and design values $\pi_T = 0.45$, $\pi_C = 0.30$.

- (a) What is the probability that the trial will fail to demonstrate a true effect?

[1]

Type II error (we were underpowered)

$$1 - 0.85 = 0.15$$

- (b) What is the probability of improvement under the standard treatment assumed to be in this trial?

[1]

$$P(\text{IMPROVE} \mid \text{control}) = \pi_c = 0.3$$

- (c) In the trial, half the patients will be allocated to the standard treatment and half to the new drug. What is the required sample size per group for this trial?

[3]

HB PAGE 13 BINT 2A
HB PAGE 27 TABLE 2

$$Z_\gamma = 1.036$$

$$\pi_0 = 0.375$$

$$Z_{(1-\alpha/2)} = Z_{0.995} = 2.576$$

$$\text{So } n = 271.8$$

So we need 272 patients in each arm of the trial.

Question 6

Several randomized controlled trials have been conducted to assess the effectiveness of propranolol (a beta blocker) at reducing mortality in patients who have suffered a heart attack. Eight studies with similar methodologies were combined in a meta-analysis and the resulting forest plot is shown in Figure 1. With each study an odds ratio of less than 1 indicates that propranolol reduced mortality.

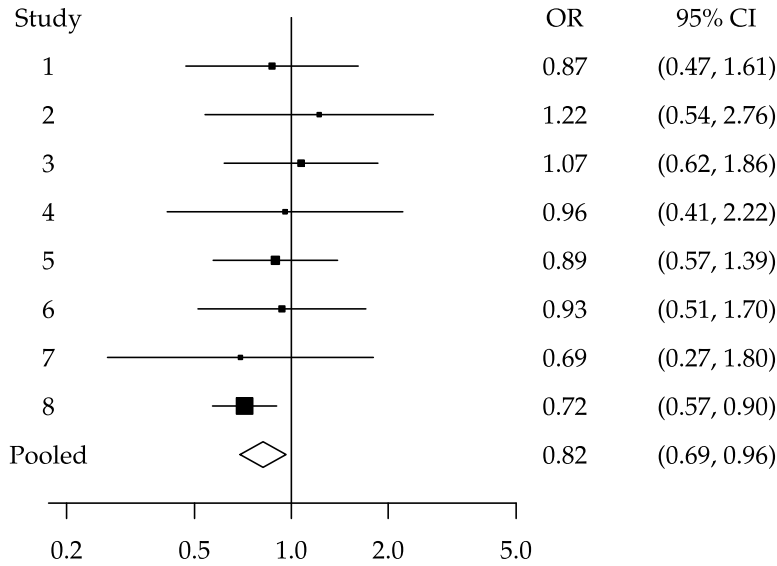


Figure 1

(a) Which study contributed most to the pooled odds ratio? Give a reason for your answer.

[1]

Study 8 (size of black square, smallest CI.)

(b) Interpret the pooled estimate of the odds ratio and its 95% confidence interval.

[1]

Pooled estimate of the CI is (0.69, 0.96) which is entirely below 1, so the result is significant at the 95% confidence level.

The OR estimate of 0.82 suggests that propandol reduces mortality.

For Examiner's use only:

Question No.	1	2	3	4	5	6	Total for Part 1
Mark							

Section 2 (relates to Book 2 *Time Series*)

Questions 7 to 14

You should **attempt all questions**. *This section is worth 25%*.

Write your answers in the spaces provided.

Question 7

Figure 2 shows the quarterly time series plot of the number of visits abroad by UK residents between 2015 and 2018.

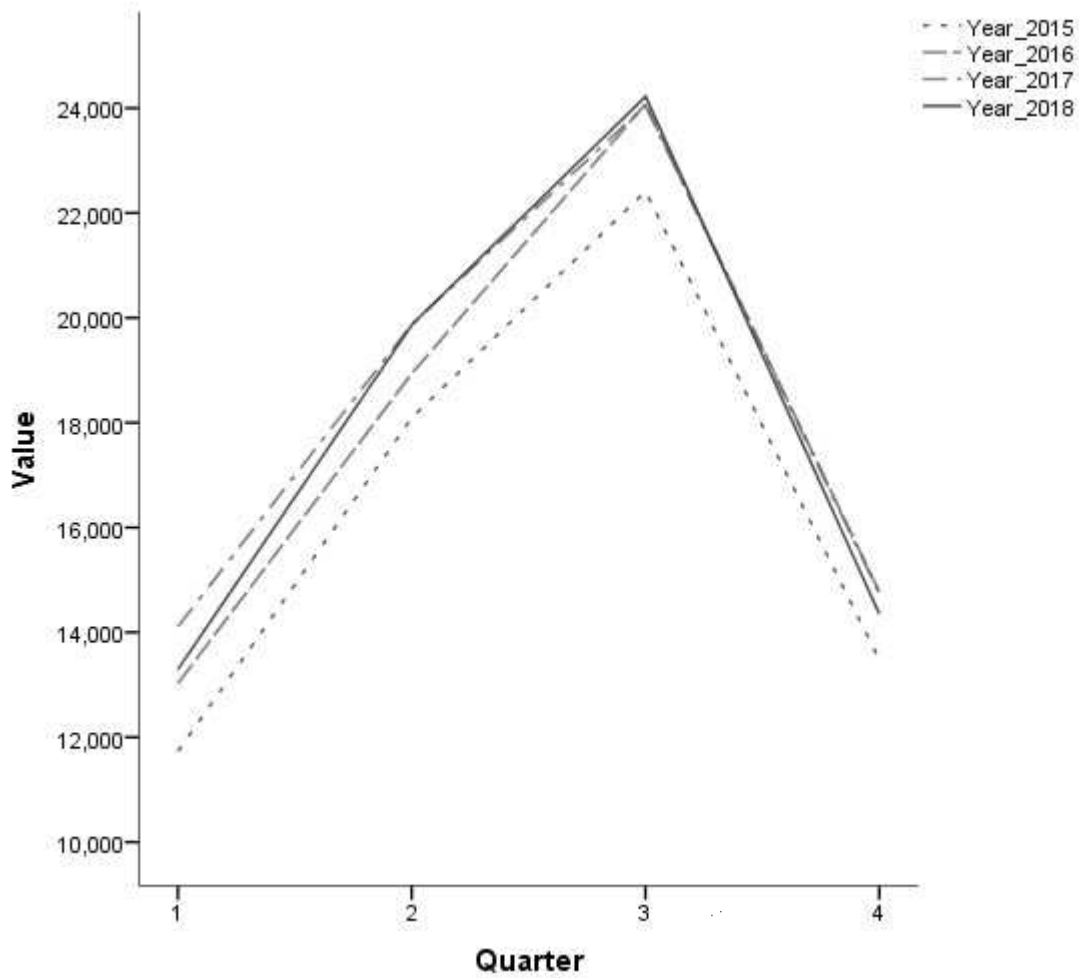


Figure 2

(a) Describe the seasonal variation in this time series.

[2]

Consistent annual pattern visits lowest in the first and fourth quarters, peaking at 24,000 in Q3.

(b) Is there any evidence of an increasing or decreasing trend from year to year in the period 2015–2018? Explain your answer.

[2]

The plot for 2015 is clearly below the other lines in each Q, with a maximum of 22,000, followed by 2016 and with an increasing trend to 2018 with the highest value (24,000).

Question 8

Figure 3 shows a time plot of the monthly UK consumer price index (CPI) inflation rate figures from January 2015 until December 2018.

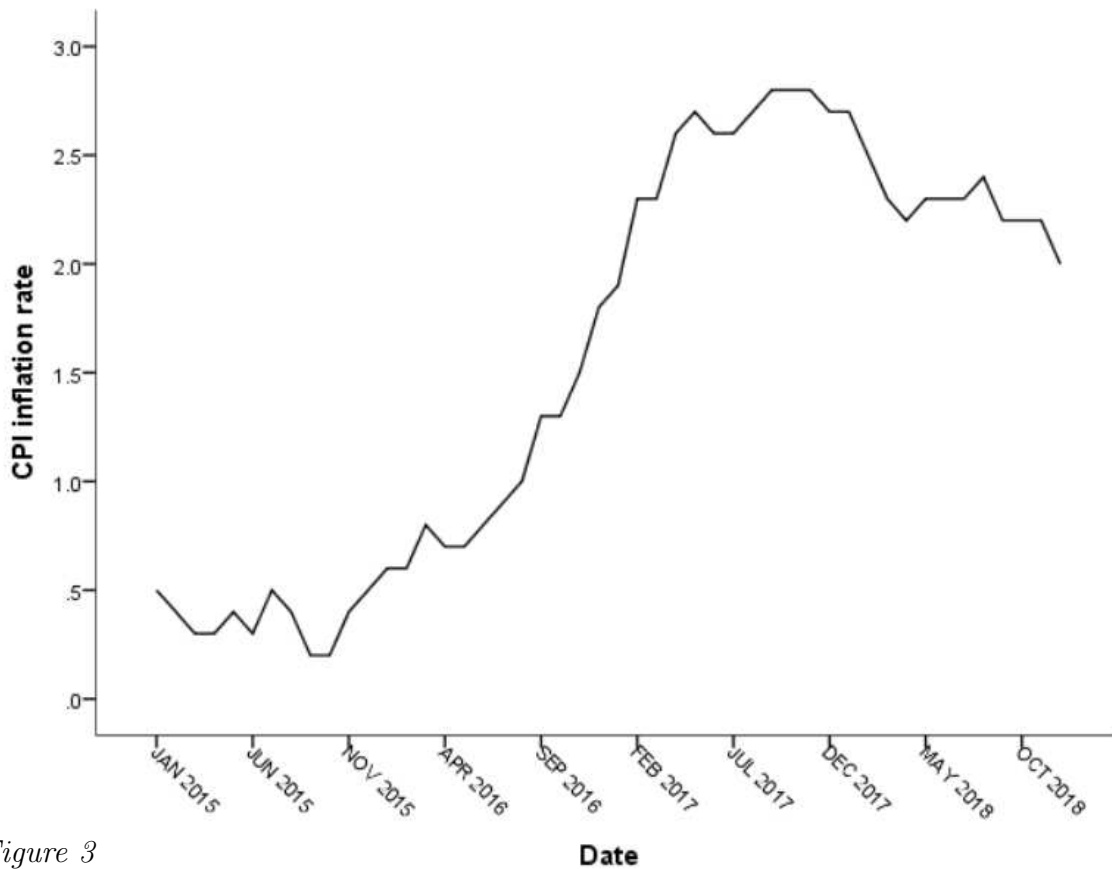


Figure 3

State which exponential smoothing model should be used for forecasting the CPI inflation rate: simple, Holt's or Holt-Winters. Explain why your chosen model is appropriate, and why each of the other two models is not appropriate.

[2]

Has PAGE 15, point 13.

Level is not constant \Rightarrow Simple not appropriate
In context CPI has seasonal variation, based on goods with seasonal variation in demand \Rightarrow Holt's not appropriate

Hence choose Holt-Winters

Question 9

Consider the quarterly seasonal time series given in Figure 4.

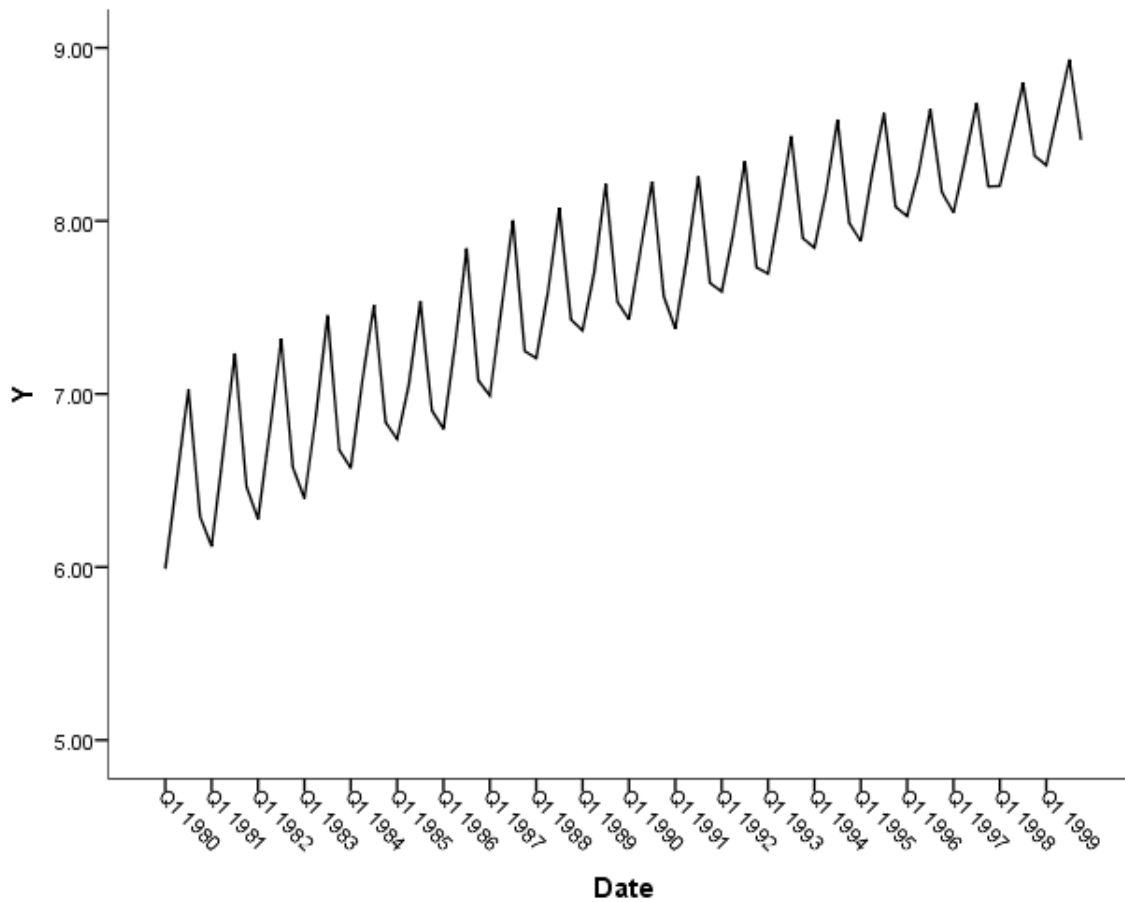


Figure 4

Which model is more appropriate for the time series in Figure 4, an additive decomposition model or a multiplicative decomposition model?

Explain your answer.

[1]

The magnitude of the fluctuations appears to be decreasing,
Hence additive model not appropriate, multiplicative model may be appropriate

Question 10

The estimated seasonal factors for the monthly series of expenditure abroad (in £ millions) by UK residents from January 2012 to December 2017 are shown in Table 4.

Table 4

Month	Seasonal factor	Month	Seasonal factor
1	-452.34	7	730.26
2	-722.09	8	?
3	-468.19	9	370.99
4	-139.44	10	42.82
5	58.64	11	-285.49
6	243.62	12	-116.94

(a) Obtain the seasonal factor for August (Month 8).

[1]

Factors sum to zero, so:

$$738.16$$

(b) The expenditure figure for August 2017 was £3014.00 millions. Calculate the seasonally adjusted expenditure figure for that month.

[1]

$$3014 - 738.16 = \pounds 2275.84 \text{ millions}$$

Question 11

The optimal values of the smoothing parameters α , γ and δ obtained with the Holt–Winters exponential smoothing method applied to a monthly time series were $\alpha = 0.32$, $\gamma = 0.71$ and $\delta = 0.01$. Interpret these optimal parameters in terms of the weight given to recent observations.

ADD TO HANDBOOK (SECTION 7.2 of BOOK 2)

[2]

Each parameter is between 0 and 1.
Values close to 1 \Rightarrow much weight is placed on recent observations.
So here γ has much weight on recent obs, less so for α & very little for δ .

Question 12

A monthly time series has 36 observed values from January 2016 to December 2018. The first two values of the log transformed series are 7.98 in January 2016, and 7.93 in February 2016.

- (a) Assuming a Holt's exponential smoothing method is used for the log transformed time series, suggest appropriate initial values for the level and the slope.

[2]

Sensible choice for the initial value for the level of the log series is the first value, i.e. 7.98

For the slope, can use $x_2 - x_1$, i.e.

$$7.93 - 7.98 = \underline{-0.05}$$

- (b) Table 5 shows the values of the SSE obtained for several pairs of values of the parameters α and γ , and the corresponding forecasts for January 2019.

Table 5

Model	α	γ	SSE	Forecast
A	0.39	0.001	2.609	7.83
B	0.39	0.005	2.605	7.81
C	0.40	0.001	2.601	7.79
D	0.40	0.005	2.602	7.80
E	0.41	0.001	2.676	7.82
F	0.41	0.005	2.763	7.85
G	0.50	0.01	2.802	7.73

- (i) Identify the optimal combination of parameter values among those listed in Table 5. Briefly explain your choice.

[2]

Equal number of parameters, so just choose minimum SSE, i.e. MODEL C (2.601)

- (ii) Obtain the forecasted value of the original series according to the optimal Holt's model for January 2019, rounded to two decimal places.

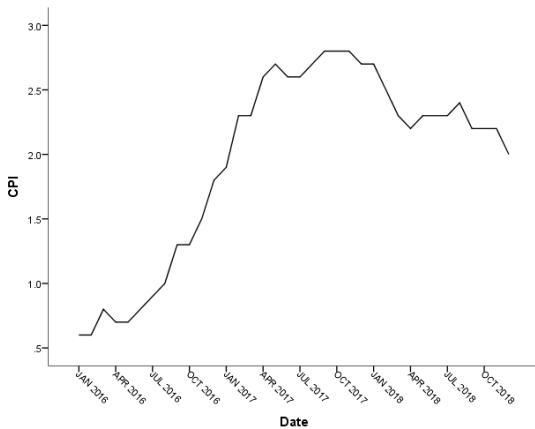
[2]

So the log-transformed prediction is 7.79.

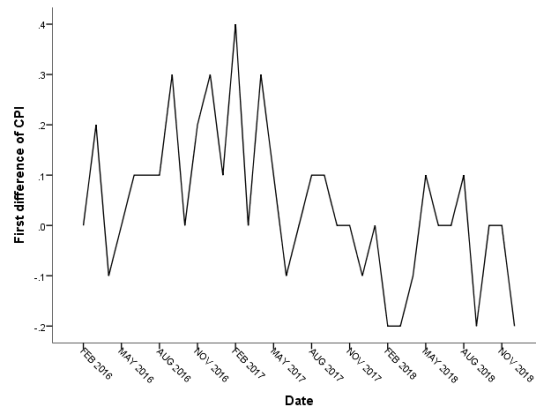
Hence the original series prediction is $\exp(7.79) = \underline{\underline{2416.32}}$

Question 13

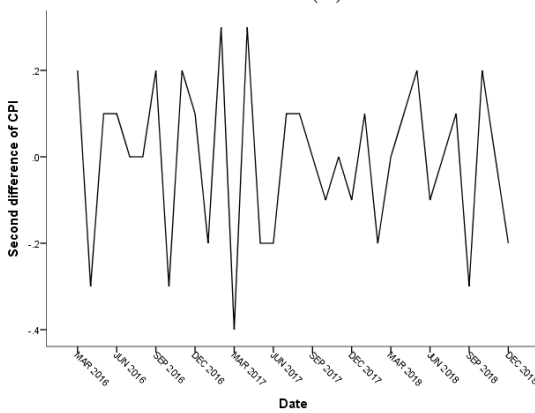
Figure 5(a) shows the time plot of the UK consumer price index (CPI) inflation rate from January 2016 to December 2018. Figures 5(b) and 5(c) show time plots of the first and second differences, respectively.



(a)



(b)



(c)

Figure 5

Identify the order of differencing required to obtain stationarity in mean and variance. Carefully explain your answer.

[2]

In mean trend in (a).

There is still a negative trend in (b) mean from Feb 2017 to Nov 2018 so $d=1$ is insufficient.

(c) is stationary in mean so further differencing is unnecessary.

$d=2$

Question 14

Figure 6 shows the correlogram (ACF) and the partial autocorrelogram (PACF) of a stationary time series at lags 1 to 20.

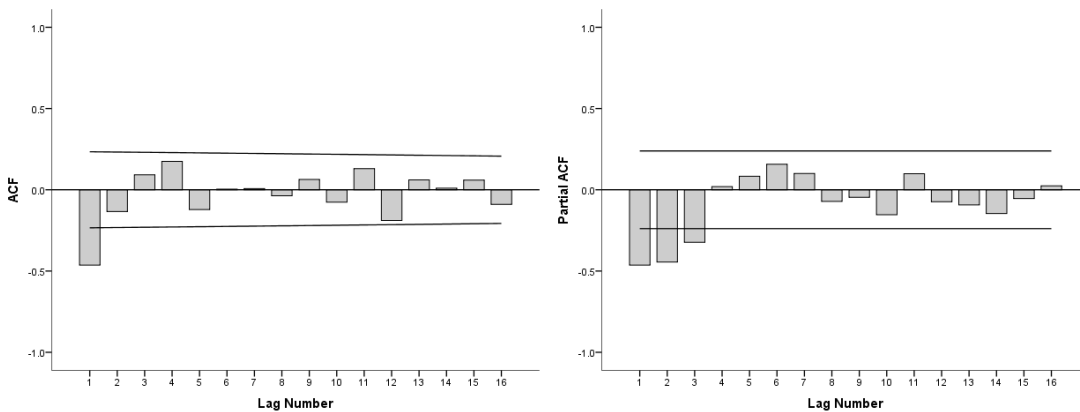


Figure 6

- (a) Identify one plausible $AR(p)$ and one plausible $MA(q)$ model for this time series. Explain carefully the reasons for your choices.

[4]

PACF cuts off after lag 3, subsequently within significance bounds suggests $AR(3)$ may be OK

ACF after lag 1 the sample autocorrelations are within significance bounds \Rightarrow $MA(1)$ may be OK

(b) Would a mixed ARMA model be a plausible model for this time series? Briefly explain your answer.

[2]

The theoretical ACF & PACF of eg. ARMA(1,1) tail off to zero with increasing lag.

This is a plausible interpretation of the correlograms, so a mixed model might be appropriate.

For Examiner's use only:

Question No.	7	8	9	10	11	12	13	14	Total for Part 2
Mark									

Section 3 (relates to Book 3 *Multivariate analysis*)

Questions 15 to 20

You should attempt all questions. *This section is worth 25%.*

Write your answers in the spaces provided.

Question 15

Table 6 shows data on:

X_1 : a baby's intake of breast milk (ml/24 hours)

X_2 : the amount of milk substitute given to the baby (ml/24 hours)

X_3 : the weight of the baby (kg)

X_4 : the weight (kg) of the baby's mother

X_5 : the height (cm) of the baby's mother.

Table 6

X_1 Milk intake (dl/24 hrs)	X_2 Supplement (dl/24 hrs)	X_3 Weight of infant (kg)	X_4 Weight of mother (kg)	X_5 Height of mother (cm)
842	250	5.002	65	173
844	0	5.128	48	158
841	40	5.445	62	160
770	0	5.667	63	165
965	60	5.106	55	162
644	240	5.196	58	170
629	0	5.526	56	153
979	30	5.928	78	175
843	0	5.263	57	170
805	230	6.578	57	168
648	555	5.588	58	173
764	0	4.613	58	171
491	0	4.360	49	162
873	60	5.882	54	163
771	315	5.618	68	179
839	370	6.032	56	175
932	130	6.030	62	168
678	0	4.727	59	172

Suppose that the data matrix based on Table 6 is denoted by \mathbf{X} .

(a) What are the values of x_{34} and x_{51} ? [1]

$$x_{34} = 5.667$$

$$x_{51} = 173$$

(b) Should the data in Table 6 be standardized before carrying out a principal component analysis? Give two reasons why, or why not. [2]

Standardization is mandatory

- different scales
- different variances

(c) The mothers' heights in the data set have a mean of 167.61cm and a standard deviation of 6.84 cm. Calculate the standardized value of mother's height for the first observation in the data set. [1]

$$\frac{173 - \bar{x}}{\sigma} \approx 0.779$$

(HB PAGE 18 POINT 6)

Question 16

The lower triangle of the sample correlation matrix from the data in Table 6 is:

$$\begin{array}{l} X_2 \\ X_3 \\ X_4 \\ X_5 \end{array} \begin{pmatrix} & X_1 & X_2 & X_3 & X_4 \\ -0.101 & & & & \\ 0.464 & 0.369 & & & \\ 0.404 & 0.125 & 0.324 & & \\ 0.143 & 0.562 & 0.140 & 0.570 & \end{pmatrix}$$

- (a) Between which pair of variables does there appear to be the strongest linear relationship. Between which pair of variables does there appear to be the weakest relationship?

[1]

max abs value is 0.57 \Rightarrow X_4, X_5 have strongest
min abs value is -0.101 \Rightarrow X_1, X_2 have weakest

- (b) Use this correlation matrix to describe the relationship between the variables.

[2]

X_5 has a moderately strong, positive linear relationship with both X_2 and X_4 .
 X_1 is positively correlated with X_3 and X_4 , as are X_2 and X_3 .
Other relationships are weak, the only negative correlation (weak) is X_1, X_2 .

Question 17

A principal component analysis was carried out on the standardized values of the variables presented in Table 6 of Question 15. The variance explained by each of the principal components (with one value missing) is given in Table 7, and the corresponding scree plot is given in Figure 7.

Table 7

Principal Component	Variance explained
1	2.230
2	1.290
3	
4	0.410
5	0.164

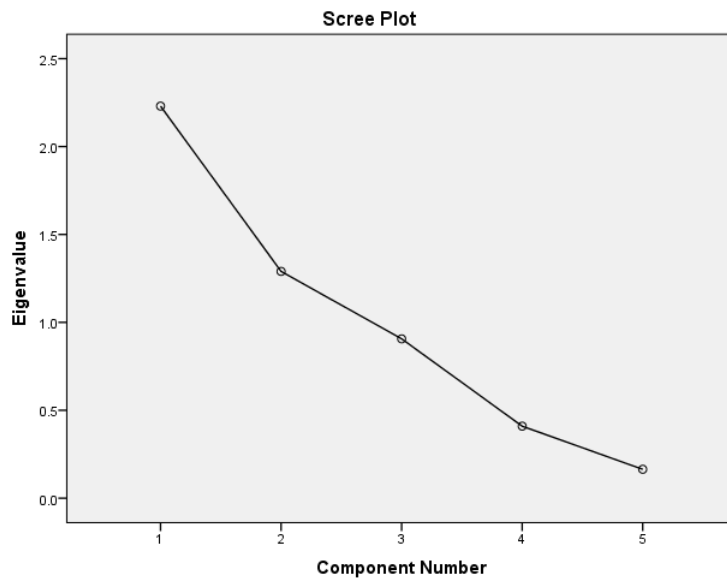


Figure 7

Use Table 6 and Figure 7 to answer the following questions.

- (a) Calculate, to 3 decimal places, the variance explained by the third principal component.

[2]

HIB PAGES 19 POINT 6

$$TV = \sum VAA = 5 \text{ (standardized)}$$
$$\text{So } VAA(PC_3) = 5 - (\sum 1,2,4,5 VAA) = 0.906$$

- (b) Using Kaiser's criterion, how many components should be retained? Explain your answer.

[2]

HB PAGE 19 POINT 15

Here average variance = 1 (standardisation)

So we retain all components with $\text{var} > 1$, i.e. $\underline{\underline{2}}$

- (c) Can you identify how many components should be retained using the scree plot. Explain your answer.

[2]

look for an 'elbow' in the scree plot, retain points before then.

Here no clear point the scree plot flattens out, so can't identify elbow.

- (d) Calculate the percentage variance explained by the first two components.

[1]

$$\text{PVE} = \frac{v(Y)}{TV} \times 100\%$$

$$= \frac{3.52}{5} \times 100\% = 70.4\%$$

Question 18

Data were collected from 50 people on the value they placed on four types of precious stone. The data were standardised and then a principal component analysis carried out using SPSS. Table 8 gives the loadings for the first two principal components produced by SPSS.

Table 8

Variable	Component 1	Component 2
Diamond	0.451	-0.631
Ruby	0.578	-0.082
Sapphire	0.537	0.121
Emerald	0.417	0.762

Briefly interpret each of the components.

[4]

Component 1:

All loadings are positive, so this is a weighted average, significant weight is given to each of the four variables.

Component 2:

Here some loadings are positive, and some negative so this is a contrast. Relatively little weight is given to ruby & sapphire, so this is essentially a contrast between diamond and emerald.

Question 19

Three variables Y_1 , Y_2 and Y_3 were recorded for a set of observations divided into four groups. The within-groups covariance matrix \mathbf{W} and the between-groups covariance matrix \mathbf{B} calculated using these data are as follows.

$$\mathbf{W} = \begin{pmatrix} 21.54 & 7.54 & 13.62 \\ 7.54 & 9.38 & 4.49 \\ 13.62 & 4.49 & 13.62 \end{pmatrix}$$

$$\mathbf{B} = \begin{pmatrix} 25.69 & -8.11 & 67.17 \\ -8.11 & 4.61 & -23.26 \\ 67.17 & -23.26 & 177.68 \end{pmatrix}$$

- (a) Calculate the separations achieved by each of the three variables on its own, and hence identify the variable that achieves the best separation between the four groups. [2]

HB PART 20 POINT 17.

$$\text{sep} = \frac{V_b}{V_w} \quad \text{so} \quad Y_1 = \frac{25.69}{21.54} = 1.19$$

$$Y_2 = 0.49$$

$$Y_3 = 13.05 \quad \leftarrow \text{best separation is } Y_3$$

- (b) Can the separation be improved substantially by using a linear combination of more than one variable? Explain your answer. [2]

Looking at the covariances, there is a different pattern of positive & negative covariances for \mathbf{W} & \mathbf{B} .

$\Rightarrow Y_i$ have different relationships between groups than within groups

\Rightarrow a linear combination of more than one variable can improve separation substantially.

Question 20

A data set consists of five variables recorded for 244 people whose work was classified as *personnel*, *mechanics* or *dispatchers*.

- (a) What is the maximum number of discriminant functions that are useful for discriminating between the three classifications using the five variables?

[1]

ADD TO HANDBOOK [BOOK 3 SECTION 12.2]

$$\min(k-1, p) = \min(2, 5) = \underline{\underline{2}}$$

- (b) Stacked histograms of the first discriminant function are shown in Figure 8 (top: *personnel*, middle: *mechanics*, bottom: *dispatchers*).

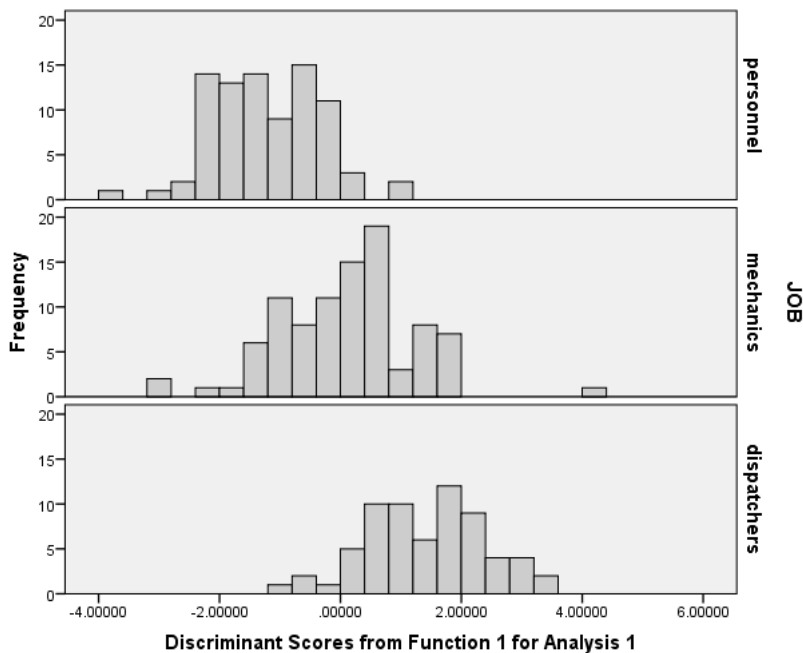


Figure 8

Briefly comment on the usefulness of this discriminant function for distinguishing between the three types of work.

[2]

The discriminant function poorly discriminates between the three types of work. This is because there is substantial overlap between the histograms, at values near 0.

For Examiner's use only:

Question No.	15	16	17	18	19	20	Total for Part 3
Mark							

Section 4 (relates to Book 4 *Bayesian statistics*)

Questions 21 to 27

You should **attempt all questions**. *This section is worth 25%.*

Write your answers in the spaces provided.

Question 21

Three different machines were used for producing a large lot of similar manufactured items. Twenty percent of the items were produced by machine A, 30% by machine B and 50% by machine C. Suppose that 3% of items produced by machine A are defective, 4% produced by machine B are defective and 2% produced by machine C are defective.

(a) What proportion of items in the complete lot are defective?

[1]

$$\begin{aligned} P(F) &= P(A)P(F|A) + P(B)P(F|B) + P(C)P(F|C) \\ &= (0.2)(0.03) + (0.3)(0.04) + (0.5)(0.02) \\ &= 0.028 \end{aligned}$$

Suppose an item is randomly selected from the complete lot.

(b) What is the probability it was produced by machine A and is a defective item?

[1]

$$P(A \cap F) = P(A)P(F|A) = (0.2)(0.03) = 0.006$$

(c) If the item is defective, what is the probability it was produced by machine A?

[2]

$$P(A|F) = \frac{P(A \cap F)}{P(F)} = \frac{0.006}{0.028} = 0.2143$$

Question 22

Figure 9 contains plots of four prior distributions.

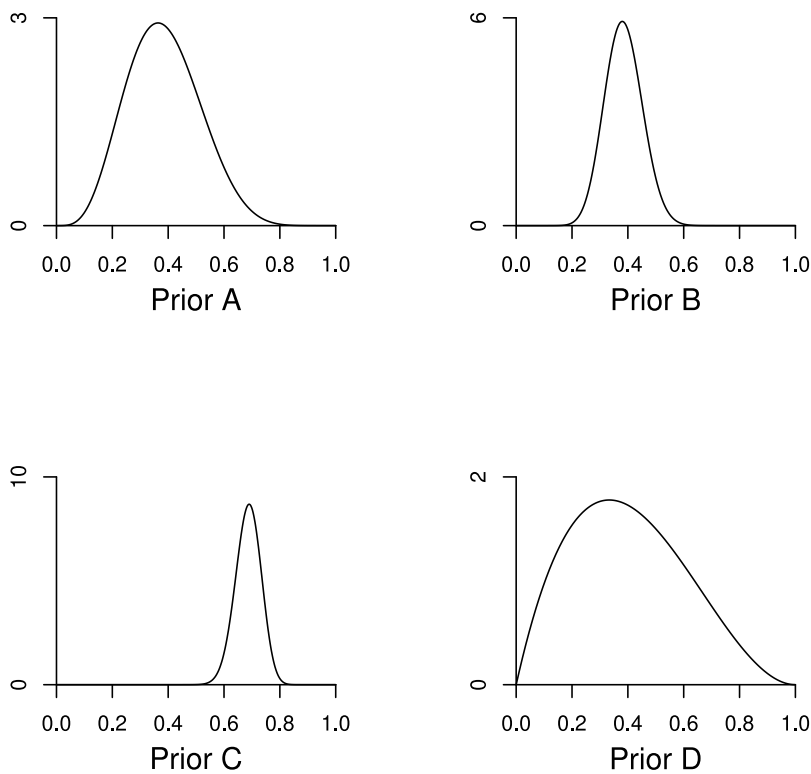


Figure 9

Place these prior distributions in order from weakest to strongest.
Justify your answer.

[2]

(weakest) D A B C

density of a weaker prior is spread over a greater range of values, so they are wider & flatter.

Question 23

Pat is about to take up a new job in which she will sometimes have to go on emergency calls. Suppose that the number of calls she will have to answer in a week follows a Poisson distribution and let θ be the average number per week.

- (a) Pat decides to represent her opinion about θ by a gamma prior distribution, $\text{Gamma}(a, b)$. Give one reason why a gamma prior is an appropriate choice. [1]

$$\theta \in \mathbb{R}^+$$

- (b) Pat chooses $a = 4$ for the first parameter of the gamma prior for θ . What value should she choose for b if she expects to make 1.5 calls per week, on average? [1]

HB PAGE 6

$$4/b = 1.5 \Rightarrow b = 2.67$$

- (c) In her first five weeks, Pat handles 2, 1, 3, 2 and 4 calls. State the posterior for θ , given the data and Pat's prior. [1]

HB PAGE 24, POINT 15

$$\begin{aligned} \text{gamma}(a + n\bar{x}, b + n) &= \text{gamma}(4 + 12, 2.67 + 5) \\ &= \text{gamma}(16, 7.67) \end{aligned}$$

- (d) Discuss how the numbers of calls in the first five weeks have changed Pat's beliefs. [2]

$$\text{If } \theta \sim \text{gamma}(16, 7.67), \quad E(\theta) = 2.09$$

So her expected number of calls has increased.

As $\text{VAR}(\theta) = a/b^2$, her confidence in this estimate has also increased.

Question 24

Suppose that data X can be modelled by the normal distribution $N(\mu, 30)$. Jim's prior for μ has median 100. The interval (97, 103) gives the central 50% of his prior.

- (a) Calculate the values of the parameters a and b of the Normal prior $N(a, b)$ that best matches Jim's beliefs about μ .

[3]

Median of normal distribution = mean, so $a = 100$.
0.75 quantile of normal distribution is 0.68σ , HB PAGE 26
So here $0.68\sigma = 3 \Rightarrow b = \sigma^2 = 19.46$
So $N(100, 19.46)$

- (b) Three observations on X are made and their sample mean is 102.0. Using the prior you specified in part (a), calculate the posterior for μ .

[3]

HB PAGE 24, POINT 15

$$a' = \frac{\sigma^2 a + n\bar{x}}{\sigma^2 + nb} = \frac{30(100) + 3(19.46)102}{30 + 3(19.46)} = 101.32$$

$$b' = \frac{\sigma^2 b}{\sigma^2 + nb} = \frac{30(19.46)}{30 + 3(19.46)} = 6.61$$

So $\mu | \text{data} \sim \text{Normal}(101.32, 6.61)$

- (c) Janet's prior for μ has mean 101, and its lower and upper quartiles are 99 and 110. Give one reason why Janet's beliefs cannot adequately be represented by a normal prior.

[1]

Not symmetric about the mean.

Question 25

Write down the model corresponding to the following WinBUGS model definition.

```
model
{
  x ~ dbin (p,n)
  p ~ dbeta (a,b)
  n <- 12
  a <- 5
  b <- 8
}
```

[2]

COMPUTER BOOK 4 ACTIVITY 8.2

$x \sim \text{Binomial}(12, p)$

where

$p \sim \text{beta}(5, 8)$

(a beta/binomial conjugate model)

Question 26

Samples from the posterior distributions of parameters α and β were obtained using WinBUGS. Using these samples, the following numerical summaries were produced.

Table 9

node	mean	sd	MC error	2.5%	median	97.5%	start	sample
alpha	-60.8	5.137	0.2208	-70.6	-60.87	-50.92	1	1000
beta	34.31	2.888	0.1252	28.72	34.35	39.93	1	1000

- (a) Write down an estimate of the posterior mean and 95% equal-tailed credible interval for β .

[1]

34.31, 95% CI (28.72, 39.93)

- (b) Describe, in general terms, what would happen to the values of sd and MC error if the number of iterations were to increase substantially. Explain your answer.

[2]

MC error $\propto \frac{1}{\sqrt{n}}$ so will decrease

as n increases.

Assuming equilibrium distribution, chain has constant mean and variance, so sd will be similar.

Question 27

Using WinBUGS, a chain of 1000 iterations was run so as to sample from the posterior for a parameter β . The resulting trace plot for β is given in Figure 10.

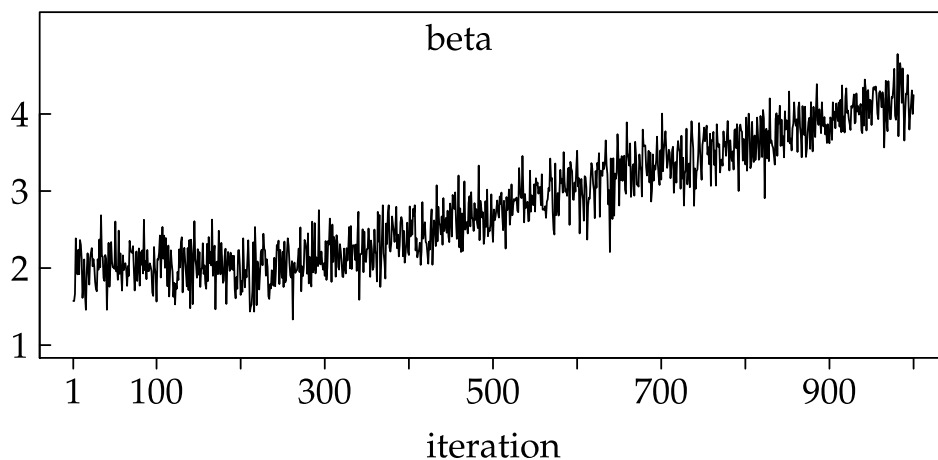


Figure 10

Do you think the chain converged? Justify your answer.

[2]

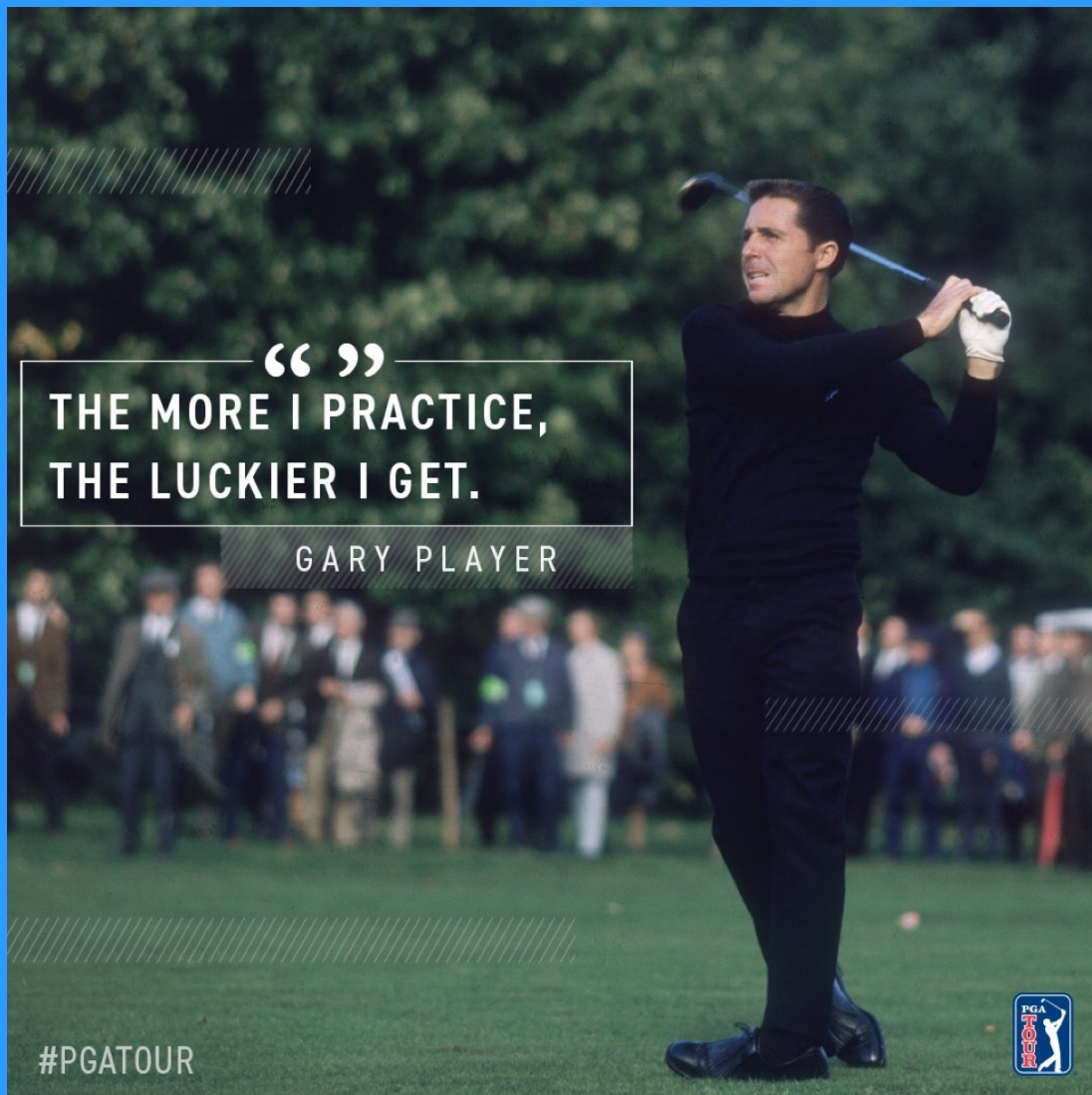
No - values do not oscillate about a common mean

For Examiner's use only:

Question No.	21	22	23	24	25	26	27	Total for Part 3
Mark								

[END OF QUESTION PAPER]

Good luck!



<https://twitter.com/pgatour/status/938800059973586944?s=12>