

A Level Mathematics Notes - Statistics

Xingzhi Lu

2129570@concordcollege.org.uk

Contents

I AS Statistics	2
1 Data collection	3
2 Measurements of location and spread	10
3 Representations of data	12
4 Correlation	13
5 Probability	14
6 Statistical distributions	15
7 Hypothesis testing	17
II A2 Statistics	20
1 Regression, correlation and hypothesis testing	21
2 Conditional probability	22
3 The normal distribution	23

Part I

AS Statistics

Chapter 1

Data collection

1.1 Populations and samples

1.1.1 Definitions

Population The whole set of items that are of interest

Census Observes or measures every member of a population

Sample A selection of observations taken from a subset of the population which is used

Sampling unit Individual units of a population that can be sampled

Sampling frame A list of all people or item that can potentially be involved in the sample

1.1.2 Census

Advantages

- Gives a completely accurate result, no bias

Disadvantages

- Time consuming and expensive
- Cannot be used when the testing process destroys the item
- Hard to process large quantity of data

1.1.3 Sample

Advantages

- Easier to implement
- Quicker to implement
- Less data to process
- Cheaper to implement

Disadvantages

- The data may not be representative
- The sample may not be large enough to give information about small sub-groups of the population

1.1.4 Sample size

- Larger sample size = better accuracy
- If the population is varied a larger sample size is needed to make sure that the sample is representative

1.2 Random sampling methods

1.2.1 Simple random sampling

Definition

- Every possible sample of size n has an **equal chance** of being picked

Method

1. Each sampling unit is numbered from 1 to n
2. Generate x random number between 1 to n using random number generators / lottery picks / random number tables (or draw out x names from the lottery hat), ignoring repeats
3. Sampling units corresponding to these numbers become the sample
4. Data taken from the sample

Advantages

- Free of bias
- Easy and cheap to implement for small populations and small samples
- Each sampling unit has a known and equal chance of selection

Disadvantages

- Not suitable when the population size or the sample size is large as it is potentially time consuming, disruptive and expensive
- A sampling frame is needed
- Chance of being unrepresentative

1.2.2 Systematic sampling

Definition

- The required elements are chosen at **regular intervals** from an **ordered list**

Method

1. The population is ordered with a unique number each from 1 to n
2. Required elements are chosen at regular intervals i.e. take every k th elements where $k = \frac{\text{Population size}}{\text{Sample size}}$
3. Starting at random item between 1 and k using a random number generator
4. Take that item and select the remaining data at the chosen interval

* **Show working**

Advantages

- Simple and quick to use
- Suitable for large samples and large populations

Disadvantages

- A sampling frame is needed
- It can introduce bias if the sampling frame is not random

1.2.3 Stratified sampling**Definition**

- The population is divided into mutually exclusive strata and a random sample is taken from each

Method

1. Population divided into **non-overlapping** groups / strata
2. Same proportion ($\frac{\text{Sample size}}{\text{Population size}}$) sampled from each strata (**show working** for the total population and the size of each strata individually, round if needed)
3. Simple random sampling carried out in each group (explain in more details here)

Advantages

- Sample accurately reflects the population structure
- Guarantees proportional representation of groups within a population

Disadvantages

- Population must be clearly classified into distinct strata
- Selection within each stratum suffers from the same disadvantages as simple random sampling

1.3 Non-random sampling methods**1.3.1 Quota sampling****Method**

1. Population divided into groups according to a given characteristic
2. A quota group is set to try and reflect the group's proportion in the whole population
3. An interviewer or researcher selects a sample that reflects the characteristics of the whole population (opportunity sampling)

* **Show working**

Advantages

- Allows a small sample to still be representative of the population
- No sampling frame required
- Quick, easy, inexpensive
- Allows for easy comparison between different groups of population

Disadvantages

- Non-random sampling can introduce bias
- Population must be divided into groups, which can be costly or inaccurate
- Increasing scope of study increases number of groups, adding time or expense
- Non-responses are not recorded

1.3.2 Opportunity / convenience / pragmatic sampling**Method**

1. Sample taken from people who are available at time of study and meet the criteria

Advantages

- Easy to carry out
- No sampling frame required
- Inexpensive

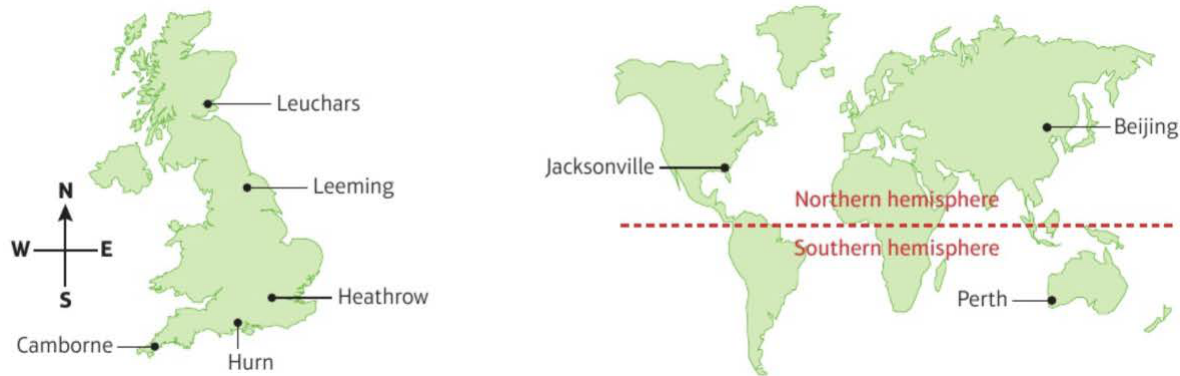
Disadvantages

- Likely to be unrepresentative
- Highly dependent on individual researcher (likely to be biased)

1.4 Large data set**1.4.1 Scope**

- Months included: May - October
- Years included: 2015 and 1987

1.4.2 Background information



City	Climate and geographical locations
Leuchars	Coastal NE of Scotland Climate generally warm and temperate significant rainfall throughout the year
Leeming	Inland Climate generally warm and temperate significant rainfall throughout the year
Heathrow	Inland Temperate oceanic climate Cool to warm summers cold winters
Hurn	Coastal Southern England Mild climate Warm summers + heavy rainfall often in mild winters
Camborne	Coastal Cornwall (SW England) Climate generally warm and temperate High rainfall even in driest months
Beijing	Inland (150km from the sea) Northern hemisphere but relatively far South, so it tends to be hot and humid in summer months
Jacksonville	Coastal Northern hemisphere but relatively far South, so it tends to be hot and humid in summer months
Perth	Coastal In the southern hemisphere - in winter during May - Oct

(Cities are ordered from North to South)

1987 “Great storm” in UK in October so there are unusually high winds, mild “El Nino” impact globally

2015 Strong “El Nino” impact espacially in the US so there is cooler temperature and higher rainfall

1.4.3 Data recorded

British locations

Variable	Unit
Daily mean temperature	The average of the hourly temperature (°C) readings, 09:00 - 09:00 GMT A reading which is not available is listed as 'n/a'.
Daily total rainfall	Daily total precipitation (mm) 09:00 - 09:00 GMT (includes snow or hail, which is melted and measured in the same way as rainfall.) 'Trace' (tr) is less than 0.05 mm. A reading which is not available will be shown by 'n/a'
Daily total sunshine	Sunshine amounts are recorded in hours and tenths and show the amount of bright sunshine recorded on the day of entry. A reading which is not available will be shown by 'n/a'
Daily maximum relative humidity	A measure of how close the air is to being saturated with water vapour. Relative humidities above 95% are associated with mist and fog. A reading which is not available will be shown by 'n/a'
Daily mean wind direction	The daily mean wind direction the wind is coming from , (clockwise from North) is averaged and rounded to the nearest 10° Readings which are not available are listed as 'n/a'.
Daily mean windspeed	Daily average windspeed Readings are taken 00:00 - 00:00 GMT, in knots (kn, 1 knot = 1.15mph) Readings which are not available are listed as 'n/a'.
Daily maximum gust	Maximum instantaneous wind speed Readings are taken 00:00 - 00:00 GMT, in knots (kn, 1 knot = 1.15mph) Readings which are not available are listed as 'n/a'.
Daily maximum gust direction	The direction from which the wind was blowing when the maximum gust during the hour commencing at the time of entry occurred, and is measured in degrees from true north. Readings which are not available are listed as 'n/a'.
Daily mean cloud cover	Measured in eights (oktas)
Daily mean visibility	The greatest distance at which an object can be seen and recognized in daylight, or at night could be seen and recognized if the general illumination were raised to daylight level. Visibility is measured horizontally, in decametres (Dm) dam = 10m A dash (-) indicates data not available.
Daily mean pressure	Mean sea level pressure, calculated from a measurement made at station level. Measured in hectopascals (hPa) where 1 hPa = 1 millibar

Overseas locations

Only the following variables are recorded:

- Daily mean temperature
- Daily total rainfall
- Daily mean windspeed (in knots and beaufort conversion)
- Daily mean pressure

1.4.4 Unit and precision of data

Variable	Unit	Precision
Daily mean temperature	°C	to 1 dp
Daily total rainfall	mm	to 1 dp (tr = less than 0.05 mm, treat as 0)
Daily total sunshine	hours	to 1 dp
Daily maximum relative humidity	as a percentage	nearest integer
Daily mean wind direction	degree + cardinal direction	nearest integer
Daily mean windspeed	knots / Beaufort conversion	nearest integer
Daily maximum gust	knots	nearest integer
Daily maximum gust direction	degree + cardinal direction	nearest integer
Daily mean cloud cover	oktas	integer from 0-8
Daily mean visibility	decametres (Dm)	nearest 100
Daily mean pressure	hectopascals (hPa)	nearest integer

1.4.5 Typical values

Temperature and wind speed

Location	Temperature range (°C)	Wind speed range (knots)
Leuchars	4-19	3-23
Leeming	4-23	3-17
Heathrow	8-29	3-19
Hurn	6-24	2-19
Camborne	10-20	3-18
Beijing	8-33	2-9
Jacksonville	15-31	1-12
Perth	8-25	4-14

Other data

Variable	Typical values
Gust	20 kn
Rainfall	0-60 mm in the UK, more extreme maximums elsewhere (e.g. 102mm in Perth)
Pressure	1013 ± 25 Pa
Wind speed on Beaufort Scale	Mostly light / moderate. Maximum is fresh (5)
Sunshine	0-16 hours
Cloud cover	0-8 oktas

1.4.6 Cleaning data

tr Needs to be replaced with a number between 0 and 0.05 (ideally 0.025 as it is the midpoint) before processing data

n/a Problem = data isn't available, usually ignored when doing calculations

Chapter 2

Measurements of location and spread

2.1 Types of means

Mean $\bar{x} = \frac{\sum x}{n}$

Median The middle value when the data values are put in order

Mode / modal class The value or class that occurs most often

2.2 Quartiles and percentiles

2.2.1 Finding medians

Ungrouped / know all individual values $\left(\frac{n+1}{2}\right)$ th value

Categorical $\frac{n}{2}$ th value

2.2.2 Quartiles

- Calculate $k = \frac{n}{4}(Q_1)$ or $k = \frac{3n}{4}(Q_3)$
- If it is an integer then the answer is $\frac{k\text{th value} + (k+1)\text{th value}}{2}$
- Otherwise take the $\lceil k \rceil$ th value

2.2.3 Percentiles

- No rounding needed, use linear interpolation straightaway

2.3 Types of data

Quantitative data Associated with numerical observations

Qualitative data Associated with non-numerical observations

Continuous data Can take any value in a given range

Discrete data Can only take specific values in a given range

2.4 Standard deviation / variance

Variance $\text{Var}(x) = \frac{\sum x^2}{n} - \left(\frac{\sum x}{n}\right)^2$

Standard deviation $\sigma = \sqrt{\text{Var}(x)} = \sqrt{\frac{\sum x^2}{n} - \left(\frac{\sum x}{n}\right)^2}$

2.5 Grouped data

2.5.1 Assumptions for estimating mean and standard deviation

- Values are **evenly distributed** within the classes

2.6 Interpreting distributions

Measuring location Mean / median / mode

Measuring spread of data Variance / standard deviation / range / interpercentile ranges

Chapter 3

Representations of data

3.1 Outliers

- An extreme value that lies outside the overall pattern of data
- By default use $LB = Q_1 - 1.5 \times (Q_3 - Q_1)$, $UB = Q_3 + 1.5 \times (Q_3 - Q_1)$

3.2 Cleaning data

- Anomalies (**not all outliers**) should be removed
- Anomalies = when the outlier is clearly an error and will be misleading

3.3 Histogram

3.3.1 Reasons for using histograms

- Data is continuous
- Data is in groups (with uneven widths)

3.3.2 Characteristics of histograms

- $\text{area} \propto \text{frequency}$

3.4 Comparing data sets

Comparing location A has a higher median / mean than B on average so A is ... than B on average

Comparing spread A has a higher IQR / standard deviation than B so there is more variation in the ... of A than B

Chapter 4

Correlation

4.1 Definitions

Bivariate data Data which has pairs of related values

Independent / explanatory variable The variable that the researcher can control, usually plotted on the x-axis

Dependent / response variable The variable that the researcher measures, usually plotted on the y-axis

Correlation Describes the nature of the linear relationship between 2 variables

4.2 Causal relationships

- 2 variables have a casual relationship if a change in 1 variable causes a change in the other
- ★ Correlation doesn't mean causation (add some explanations in context for questions)

4.3 Linear regression

* Work these out using a **calculator** in exams

4.3.1 Regression equation for least squares regression line

- Regression line of y on x : $y = a + bx$
- $b = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2}$ (**not needed for the exam**)
- $a = \bar{y} - b\bar{x}$
- Positive correlation: b positive, negative correlation: b negative

4.3.2 Predicting values

- Should not extrapolate, only do interpolation
- Reliability: reliable as it is within the range of data / not reliable as it is extrapolating
- ★ Not suitable for predicting x based on y (the independent variable in this model is x , you should not use this model to predict the value of x based on y)

4.3.3 Reason for using a regression line

- The data shows a strong (positive / negative) linear correlation

Chapter 5

Probability

5.1 Definitions

Experiment A repeatable process that gives rise to a number of outcomes

Event A set of one or more of these outcomes

Sample space A set of all the possible outcomes

Mutually exclusive Events cannot happen at the same time

Independent events Whether one event happens does not affect the probability of the other happening

Chapter 6

Statistical distributions

6.1 Definitions

Random variable One whose value depends on the outcome of a random event (outcome not known until the event took place)

Sample space The range of values that the outcome can take

Discrete variable Can only take certain numerical values

Probability distribution Fully describes the probability of any outcome in the sample space

Uniform discrete distribution All the probabilities are equal

6.2 Notations

- Capital letters (X or Y) denotes random variables
- Equivalent lowercase letters (x or y) denotes particular values of the random variable

6.2.1 Probability mass / density function

- $P(X = x) = \dots (x = \dots)$

- You might need a large bracket e.g. $P(X = x) = \begin{cases} 0.1 & x = 1, 2 \\ 0.4 & x = 3, 4 \\ 0 & x = \text{anything else} \end{cases}$

6.3 Binomial distribution

6.3.1 Notation

$$X \sim B(n, p)$$

6.3.2 Probability calculation

$$P(x) = \binom{n}{x} p^x q^{n-x}$$

6.3.3 Assumptions

- There are a fixed number of trials, n
- There are two possible outcomes only (success and failure)

- There is a fixed probability of success, p
- The trials are independent of each other

Chapter 7

Hypothesis testing

7.1 Definitions

Hypothesis A statement about the value of a population parameter

Test statistic A value computed from sample data

Null hypothesis (H_0) The hypothesis assumed to be correct ($\theta = \theta_0$)

Alternative hypothesis (H_1) Tells you about the parameter if H_0 is rejected as a result of the test ($\theta \neq \theta_0$ / $\theta > \theta_0$ (right tail) / $\theta < \theta_0$ (left tail))

Significance level (α) Probability of rejecting H_0 when assuming H_0 is true

Critical region A region of the probability distribution which, if the test statistic falls within it, would cause you to reject the null hypothesis

Critical value The first value to fall inside the critical region / a value that is compared to the test statistic to determine whether to reject H_0

Acceptance region The rejection region for H_1 in the testing of a hypothesis

Actual significance level The probability of incorrectly rejecting the null hypothesis (when H_0 is actually true)

7.2 Test on proportion / probability of success assuming binomial distribution

t^* = test statistics

7.2.1 By critical value

- One tailed: if stats test $t^* > cv$ or $t^* < cv$ (depends on right / left tail): reject H_0 , else accept H_0
- Two tailed: if stats test $t^* > \text{upper cv}$ or $t^* < \text{lower cv}$: reject H_0 , else accept H_0 (For 2 tailed tests the probability used for calculating cv at the end of each tail = $\frac{\alpha}{2}$)

7.2.2 By p value

- One tailed: if $P(t \geq t^*) < \alpha$: reject H_0 , else accept H_0
- Two tailed: if $P(t \geq t^*) < \frac{\alpha}{2}$ or $P(t \leq t^*) > \frac{\alpha}{2}$: reject H_0 , else accept H_0

7.3 Two tailed tests

- Halve the significance value to find out the critical region at each end unless otherwise specified
- Notice if the question asks for the probability in each tail to be **as close to $\frac{\alpha}{2}$ as possible**
- Always use 2 tailed tests if whether testing for increase / decrease in p is not specified

7.4 Example responses

7.4.1 One tailed + critical region

Example 7.1

A single observation is taken from $X \sim B(10, p)$ and $x = 1$ is obtained. Use this value to test $H_0 : p = 0.4$ against $H_1 : p < 0.4$ using a 5% significance level

Solution

$$H_0 : p = 0.4$$

$$H_1 : p < 0.4$$

Test statistic: $x = 1$

Significance level = 5%

One-tailed test

$$P(X \leq c_1) < 0.05$$

$$P(X \leq 1) = 0.0463 \quad (P(X \leq 2) = 0.1672 \text{ too big})$$

$$c_1 = 1 \text{ so critical region is } X \leq 1$$

$x = 1$ lies in the critical region, so evidence suggests rejecting H_0 at 5% significance level

7.4.2 One tailed + p value

Example 7.2

A single observation is taken from $X \sim B(10, p)$ and $x = 5$ is obtained. Use this value to test $H_0 : p = 0.25$ against $H_1 : p > 0.25$ using a 5% significance level

Solution

$$H_0 : p = 0.25$$

$$H_1 : p > 0.25$$

Test statistic: $x = 5$

Significance level = 5%

One-tailed test

$$\begin{aligned} P(X \geq 5) &= 1 - P(X \leq 4) \\ &= 0.0781 \end{aligned}$$

Compare p -value with significance level: $0.0781 > 0.05$

It is not significant so no evidence to reject H_0 at the 5% significance level

7.4.3 Two tailed - find critical region when using 'probability as close to'

Example 7.3

$Y \sim B(25, p)$, given that $H_0 : p = 0.42$, $H_1 : p \neq 0.42$, find the critical region for the test using 10% significance level, **the probability in each tail should be as close to 5% as possible**

Solution

$P(Y \leq c_1)$ as close to 0.05 as possible

$P(Y \leq 6) = 0.0495 < 0.05$ - closest to 0.05 so $c_1 = 6$

$P(Y \leq 7) = 0.1106 > 0.05$

$P(Y \geq c_2)$ as close to 0.05 as possible $\rightarrow 1 - P(Y \leq c_2 - 1)$ as close to 0.05 as possible $\rightarrow P(Y \leq c_2 - 1)$ as close to 0.95 as possible

$P(Y \leq 14) = 0.9465 < 0.95$ - closest to 0.05 so $c_2 = 14 + 1 = 15$

$P(Y \leq 15) = 0.19779 > 0.95$

Part II

A2 Statistics

Chapter 1

Regression, correlation and hypothesis testing

1.1 PMCC

- Measures the strength of **linear** correlation
- $r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$ (not needed for the exam)

1.2 Converting to linear form

1.2.1 Commenting on appropriateness

- ... gives a linear relationship between ...
- PMCC is close to 1 or -1 which supports the use of a linear model

1.3 Hypothesis test for zero linear correlation

1.3.1 Template

- $H_0: \rho = 0$
- $H_1: \rho \neq 0$ (two tailed) / $\rho > 0$ (right tail) / $\rho < 0$ (left tail)
- Sample size = ...
- Significance level = ...
- The critical value of r for this test is ...
- The observed value of r is ...
- ... $<$... so the observed value of r is inside / outside the critical region
- So reject / accept H_0
- Conclusion

Chapter 2

Conditional probability

2.1 Conditional probability formula

- $P(A|B) = \frac{P(A \cap B)}{P(B)}$

2.2 Principle of inclusion-exclusion

- $P(A \cup B) = P(A) + P(B) - P(A \cap B)$

Chapter 3

The normal distribution

3.1 Notation

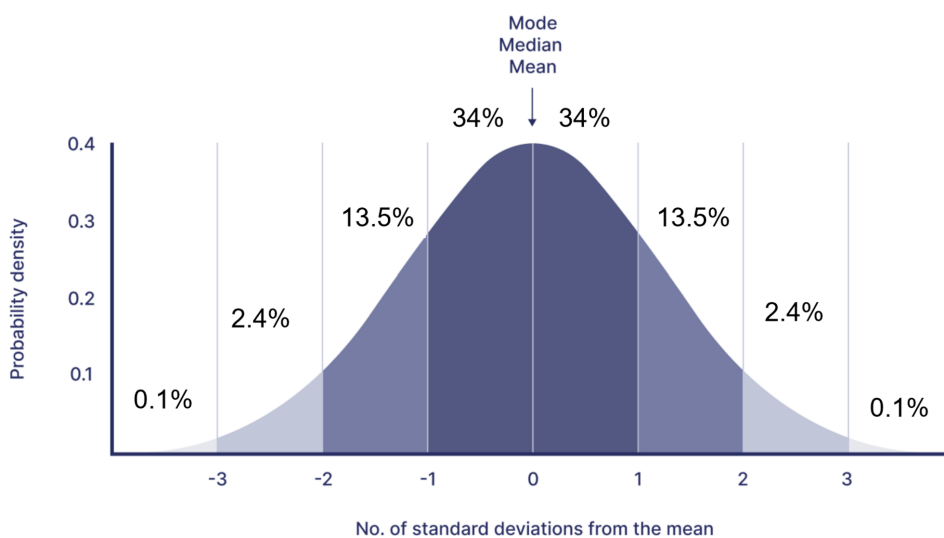
- $X \sim N(\mu, \sigma^2)$
- μ = mean of the population
- σ^2 = **variance** of the data

3.2 Properties

- The data is **continuous**
- Has parameters μ (mean) and σ^2 (variance)
- Is symmetrical: mean = median = mode
- Has a bell-shaped curve with asymptotes at each end
- Total area under the curve = 1
- Has points of inflection at $\mu + \sigma$ and $\mu - \sigma$

3.3 Estimating probabilities

- 68% of observations lie within ± 1 standard deviation of the mean
- 95% of observations lie within ± 2 standard deviation of the mean
- 99.8% of observations lie within ± 3 standard deviation of the mean



3.4 Approximation of binomial distribution

If n is large ($n \geq 35$) and p is close to 0.5, then $X \sim B(n, p)$ can be modelled as

$$Y \sim N(np, np(1-p))$$

3.4.1 Approximations

- $P(X \geq a) \approx P(Y \geq [a - 0.5])$
- $P(X = a) \approx P([a - 0.5] < Y < [a + 0.5])$
- $P(X \leq a) \approx P(Y \leq [a + 0.5])$

3.5 Sample mean

If n is large enough ($n \geq 35$) and $X \sim N(\mu, \sigma^2)$, then sample mean \bar{X} is normally distributed:

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$