



STA 2100 PROBABILITY AND STATISTICS I

PURPOSE

At the end of the course the student should be proficient in representing data graphically and handling summary statistics, simple correlation and best fitting line, and handling probability and probability distributions including expectation and variance of a discrete random variable.

DESCRIPTION

Classical and axiomatic approaches to probability. Compound and conditional probability, including Bayes' theorem. Concept of discrete random variable: expectation and variance. Data: sources, collection, classification and processing. Frequency distributions and graphical representation of data, including bar diagrams, histograms and stem-and-leaf diagrams. Measures of central tendency and dispersion. Skewness and kurtosis. Correlation. Fitting data to a best straight line.

Pre-Requisites: STA 2104 Calculus for statistics I, SMA 2104 Mathematics for Science.

COURSE TEXT BOOKS

1. Uppal, S. M., Odhiambo, R. O. & Humphreys, H. M. *Introduction to Probability and Statistics*. JKUAT Press, 2005. ISBN 9966-923-95-0
2. J Crawshaw & J Chambers *A concise course in A-Level statistics, with worked examples*, 3rd ed. Stanley Thornes, 1994 ISBN 0-534- 42362-0.

COURSE JOURNALS

1. Journal of Applied Statistics (J. Appl. Stat.) [0266-4763; 1360-0532]
2. Statistics (Statistics) [0233-1888]

FURTHER REFERENCE TEXT BOOKS AND JOURNALS

1. GM Clarke & D Cooke *A Basic Course in Statistics*. 5th ed. Arnold, 2004 ISBN13: 978-0-340-81406-2 ISBN10: 0-340-81406-3.
2. S Ross *A first course in Probability* 4th ed. Prentice Hall, 1994 ISBN-10: 0131856626 ISBN-13: 9780131856622.
3. P.S. Mann. *Introductory Statistics*. John Wiley & Sons Ltd, 2001 ISBN 13: 9780471395119.
4. Statistical Science (Stat. Sci.) [0883-4237]
5. Journal of Mathematical Sciences
6. Journal of Teaching Statistics

Introduction

What is statistics?

The Word statistics has been derived from Latin word “**Status**” or the Italian word “**Statista**”, the meaning of these words is “**Political State**” or a Government. Early applications of statistical thinking revolved around the needs of states to base policy on demographic and economic data.

Definition

Statistics: *a branch of science that deals with collection presentation, analysis, and interpretation of data.* The definition points out 4 key aspects of statistics namely

- | | |
|-------------------------|--------------------------|
| (i) Data collection | (iii) Data analysis, and |
| (ii) Data presentation, | (iv) Data interpretation |

Statistics is divided into 2 broad categories namely descriptive and inferential statistics.

Descriptive Statistics: summary values and presentations which gives some information about the data Eg the mean height of a 1st year student in JKUAT is 170cm. 170cm is a statistics which describes the central point of the heights data.

Inferential Statistics: summary values calculated from the sample in order to make conclusions about the target population.

Types of Variables

Qualitative Variables: Variables whose values fall into groups or categories. They are called categorical variables and are further divided into 2 classes namely nominal and ordinal variables

- a) Nominal variables: variables whose categories are just names with no natural ordering. Eg gender marital status, skin colour, district of birth etc
- b) Ordinal variables: variables whose categories have a natural ordering. Eg education level, performance category, degree classifications etc

Quantitative Variables: these are numeric variables and are further divided into 2 classes namely discrete and continuous variables

- a) Discrete variables: can only assume certain values and there are gaps between them. Eg the number of calls one makes in a day, the number of vehicles passing through a certain point etc
- b) Continuous variables: can assume any value in a specified range. Eg length of a telephone call, height of a 1st year student in JKUAT etc

1. Data Collection:

1.1 Sources of Data

There are 2 sources for data collection namely Primary, and Secondary data

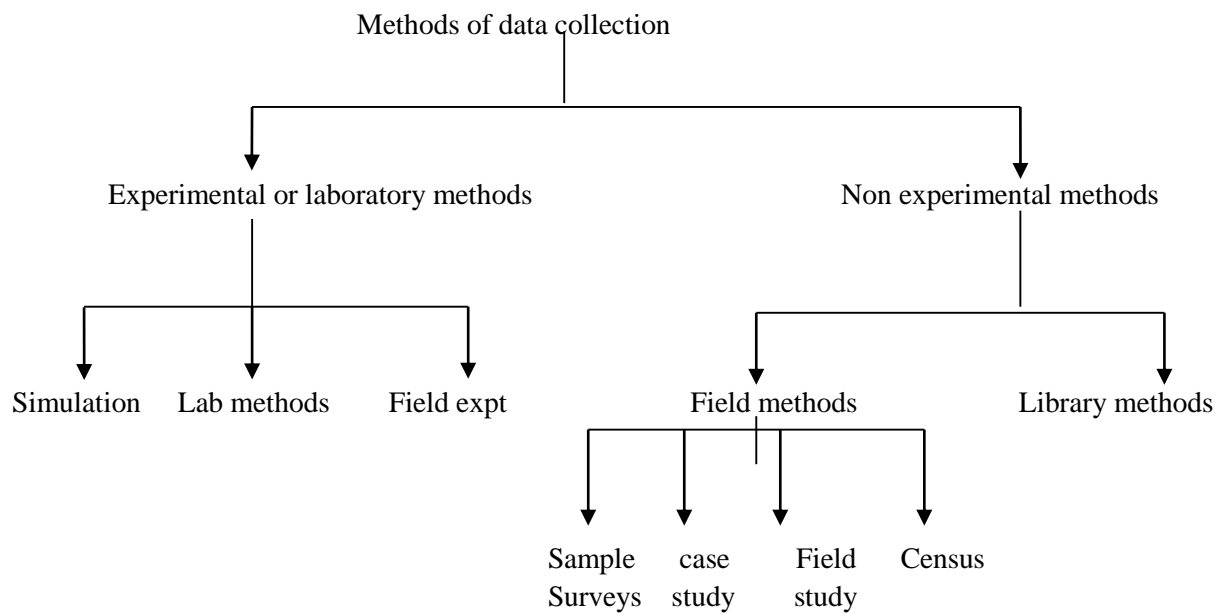
Primary data:- freshly collected ie for the first time. They are original in character ie they are the first hand information collected, compiled and published for some purpose. They haven't undergone any statistical treatment

Secondary Data:- 2nd hand information mainly obtained from published sources such as statistical abstracts books encyclopaedias periodicals, media reports eg census report CD-roms and other electronic devices, internet. They are not original in character and have undergone some statistical treatment at least once.

1.2 Data Collection Methods

The 1st step in any investigation (inquiry) is data collection. Information can either be collected directly or indirectly from the entire population or a sample.

There are many methods of collecting data which includes the ones illustrated in the flow chart below



Experimental methods are so called because in them the investigator in a laboratory tests the hypothesis about the cause and effect relationship by manipulating the independent variables under controlled conditions.

Non-Experimental methods are so called because in them the investigator does not control or change any aspect of the situation under study but simply describes what naturally occurs at a certain point or period of time.

Non-Experimental methods are widely used in social sciences. Some of the Non-Experimental methods used for data collection are outlined below.

- a) **Field study**:- aims at testing hypothesis in natural life situations. It differs from field experiment in that the researcher does not control or manipulate the independent variables but both of them are carried out in natural conditions

Merits:

- (i) The method is realistic as it is carried out in natural conditions
- (ii) It's easy to obtain data with large number of variables.

Demerits

- (iii) Independent variables are not manipulated.
- (iv) Co-operation of the organization is often difficult to obtain.
- (v) Data is likely to contain unknown sampling biasness.
- (vi) The drop rate (proportion of irrelevant data) may be high in such studies.
- (vii) Measurement is not precise as in laboratory because of influence of confounding variables.

- b) **Census**. A census is a study that obtains data from every member of a population (totality of individuals /items pertaining to certain characteristics). In most studies, a census is not practical, because of the cost and/or time required.
- c) **Sample survey**. A sample survey is a study that obtains data from a subset of a population, in order to estimate population attributes/ characteristics. Surveys of human populations and institutions are common in government, health, social science and marketing research.
- d) **Case study** –It's a method of intensively exploring and analyzing the life of a single social unit be it a family, person, an institution, cultural group or even an entire community. In this method no attempt is made to exercise experimental or statistical control and phenomena related to the unit are studied in natural. The researcher has several discretion in gathering information from a variety of sources such as diaries, letters, autobiographies, records in office, files or personal interviews.

Merits:

- (i) The method is less expensive than other methods.
- (ii) Very intensive in nature –aims at studying a few units rather than several
- (iii) Data collection is flexible since the researcher is free to approach the problem from any angle.
- (iv) Data is collected from natural settings.

Demerits

- (i) It lacks internal validity which is basic to scientific evidence.
 - (ii) Only one unit of the defined population is studied. Hence the findings of case study cannot be used as a base for generalization about a large population. They lack external validity.
 - (iii) Case studies are more time consuming than other methods.
- e) **Experiment.** An experiment is a controlled study in which the researcher attempts to understand cause-and-effect relationships. In experiments actual experiment is carried out on certain individuals / units about whom information is drawn. The study is "controlled" in the sense that the researcher controls how subjects are assigned to groups and which treatments each group receives.
- f) **Observational study.** Like experiments, observational studies attempt to understand cause-and-effect relationships. However, unlike experiments, the researcher is not able to control how subjects are assigned to groups and/or which treatments each group receives. Under this method information, is sought by direct observation by the investigator.

1.3 Population and Sample

Population: The entire set of individuals about which findings of a survey refer to.

Sample: A subset of population selected for a study.

Sample Design: The scheme by which items are chosen for the sample.

Sample unit: The element of the sample selected from the population.

Unit of analysis: Unit at which analysis will be done for inferring about the population. Consider that you want to examine the effect of health care facilities in a community on prenatal care. What is the unit of analysis: health facility or the individual woman?.

Sampling Frames

For probability sampling, we must have a list of all the individuals (units) in the population. This list or sampling frame is the basis for the selection process of the sample. "A [sampling] frame is a clear and concise description of the population under study, by virtue of which the population units can be identified unambiguously and contacted, if desired, for the purpose of the survey" - Hedayet and Sinha, 1991

Based on the sampling frame, the sampling design could also be classified as:

Individual Surveys if List of individuals is available or when the size of population is small
Special population

Household Surveys; If it's Based on the census of the households and if the individual level information is unlikely to be available In practice, it's limited to small geographical areas and known as "area sampling frame" Example: Demographic and Health Surveys (DHS)

Institutional Surveys If it's Based on the census of say Hospital/clinic lists eg

- i) 1990 National Hospital Discharge Survey
- ii) National Ambulatory Medical Care Survey

Problems of Sampling Frame

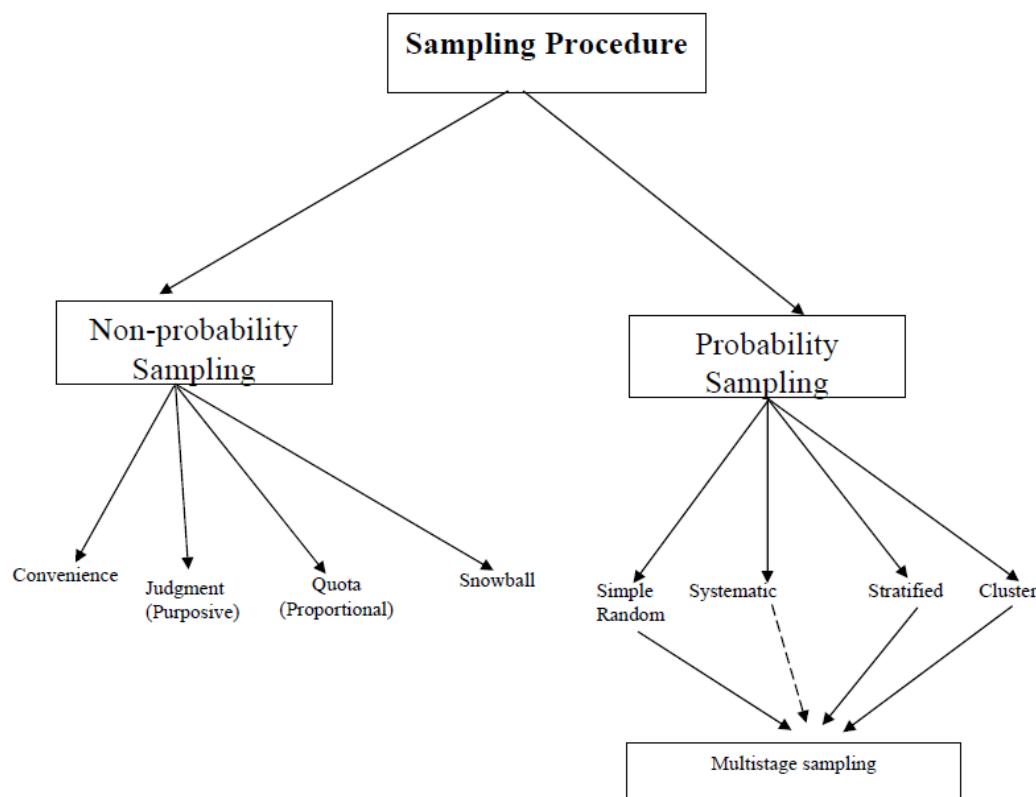
- (i) Missing elements
- (ii) Noncoverage
- (iii) Incomplete frame
- (iv) Old list
- (v) Undercoverage
- (vi) May not be readily available
- (vii) Expensive to gather

1.4 Sampling

Sampling is a statistical process of selecting a representative sample. We have probability sampling and non-probability sampling. **Probability Samples** involves a mathematical chance of selecting the respondent. Every unit in the population has a chance, greater than zero, of being selected in the sample. Thus producing unbiased estimates. They include;

- (i) Simple random sampling
- (ii) Systematic sampling
- (iii) Stratified sampling
- (iv) Cluster sampling
- (v) multi-stage sampling

Non-probability sampling is any sampling method where some elements of the population have *no* chance of selection (also referred to as “out of coverage”/“undercovered”), or where the probability of selection can't be accurately determined. It yields a non-random sample therefore making it difficult to extrapolate from the sample to the population. They include; Judgement sample, purposive sample, convenience sample: **subjective** Snow-ball sampling: **rare group/disease study**



1.4.1 Sampling Procedure

Sampling involves two tasks

- How to select the elements?
- How to estimate the population characteristics – from the sampling units?

We employ some *randomization* process for sample selection so that there is no preferential treatment in selection which may introduce selectivity bias

1.4.2 Reasons Behind sampling

- (i) Cost; the sample can furnish data of sufficient accuracy at much lower cost.
- (ii) Time; the sample provides information faster than census thus ensuring timely decision making.
- (iii) Accuracy; it is easier to control data collection errors in a sample survey as opposed to census.
- (iv) Risky or destructive test call for sample survey not census eg testing a new drug.

1.4.3 Probability Sampling Techniques

a)...Simple Random Sampling (SRS)

In this design, each element has an equal probability of being selected from a list of all population units (sample of n from N population). Though it's attractive for its simplicity, the design is not usually used in the sample survey in practice for several reasons:

- (i) Lack of listing frame: the method requires that a list of population elements be available, which is not the case for many populations.
- (ii) Problem of small area estimation or domain analysis: For a small sample from a large population, all the areas may not have enough sample size for making small area estimation or for *domain* analysis by variables of interest.
- (iii) Not cost effective: SRS requires covering of whole population which may reside in a large geographic area; interviewing few samples spread sparsely over a large area would be very costly.

Implementation of SRS sampling:

- (i) Listing (sampling) Frame
- (ii) Random number table (from published table or computer generated)
- (iii) Selection of sample

Computer generated random numbers: (STATA output)

```
832645 573158 467460 838921 171721 152885
708009 285644 727733 343305 539264 907568
305761 995036 740619 054728 746425 713746
536405 504168 750032 367682 626278 855480
217862 782003 409660 155199 129514 484511
844905 296231 103727 053603 562252 219726
670523 707073 049209 830572 337034 716264
334920 023934 808901 740693 170372 095017
885588 384435 129958 303040 264636 858065
458268 058670 888935 064613 661404 411861
277649 076177 482951 876389 898190 927367
977683 759956 553916 983998 331578 981306
```

b)..Systematic Sampling

Systematic sampling, either by itself or in combination with some other method, may be the most widely used method of sampling.” In systematic sampling we select samples “evenly” from the list (sampling frame): First, let us consider that we are dividing the list evenly into some “blocks”. Then, we select a sample element from each block.

In systematic sampling, only the first unit is selected at random, the rest being selected according to a predetermined pattern. To select a systematic sample of n units, the first unit is selected with a random start r from 1 to k sample, where $k=N/n$ sample intervals, and after the selection of first sample, every k^{th} unit is included where $1 \leq r \leq k$.

An example:

Let $N=100$, $n=10$, then $k=100/10$. Then the random start r is selected between 1 and 10 (say, $r=7$). So, the sample will be selected from the population with serial indexes of: 7, 17, 27,, 97. i.e., r , $r+k$, $r+2k$,, $r+(n-1)k$

What could be done if $k=N/n$ is not an integer?

Selection of systematic sampling when sampling interval (k) is not an integer

Consider, $n=175$ and $N=1000$. So, $k=1000/175 = 5.71$

One of the solution is to make k rounded to an integer, i.e., $k=5$ or $k=6$. Now, if $k=5$, then $n=1000/5=200$; or, If $k=6$, then $n=1000/6 = 166.67 \sim 167$. Which n should be chosen?

Solution

if $k=5$ is considered, stop the selection of samples when $n=175$ achieved.

if $k=6$ is considered, treat the sampling frame as a circular list and continue the selection of samples from the beginning of the list after exhausting the list during the first cycle.

An alternative procedure is to keep k non-integer and continue the sample selection as follows: Let us consider, $k=5.71$, and $r=4$. So, the first sample is 4th in the list. The second $= (4+5.71) = 9.71 \sim 9$ th in the list, the third $= (4+2*5.71) = 15.42 \sim 15$ th in the list, and so on. (The last sample is: $4+5.71*(175-1) = 997.54 \sim 997$ th in the list). Note that, k is switching between 5 and 6

Advantages:

Systematic sampling has many attractiveness:

- (i) Provides a better random distribution than SRS
- (ii) Simple to implement
- (iii) May be started without a complete listing frame (say, interview of every 9th patient coming to a clinic).
- (iv) With ordered list, the variance may be smaller than SRS (see below for exceptions)

Disadvantages:

- (i) Periodicity (cyclic variation)
- (ii) linear trend

i.

When to use systematic sampling?

- i) Even preferred over SRS
- ii) When no list of population exists
- iii) When the list is roughly of random order
- iv) Small area/population

c)..Stratified Sampling

In stratified sampling the population is partitioned into groups, called *strata*, and sampling is performed separately within each *stratum*.

This sampling technique is used when;

- i) Population groups may have different values for the responses of interest.
- ii) we want to improve our estimation for each group separately.
- iii) To ensure adequate sample size for each group.

In stratified sampling designs:

- i) Stratum variables are mutually exclusive (no over lapping), e.g., urban/rural areas, economic categories, geographic regions, race, sex, etc. The principal objective of stratification is to reduce sampling errors.
- ii) The population (elements) should be *homogenous* within-stratum, and the population (elements) should be *heterogeneous* between the strata.

Advantages

- (i) Provides opportunity to study the stratum; variations - estimation could be made for each stratum

- (ii) Disproportionate sample may be selected from each stratum
- (iii) The precision is likely to increase as variance may be smaller than simple random case with same sample size
- (iv) Field works can be organized using the strata (e.g., by geographical areas or regions)
- (v) Reduce survey costs.

Disadvantages

- (i) Sampling frame is needed for each stratum
- (ii) Analysis method is complex
- (iii) Correct variance estimation
- (iv) Data analysis should take sampling “weight” into account for disproportionate sampling of strata
- (v) Sample size estimation is difficult in practice

Allocation of Stratified Sampling

The major task of stratified sampling design is the appropriate allocation of samples to different strata.

Types of allocation methods:

- (i) Equal allocation
- (ii) Proportional to stratum size
- (iii) Cost based sample allocation

Equal Allocation

Divide the number of sample units n equally among the K strata. ie $n_i = \frac{n}{k}$ Example: $n = 100$ and $k = 4$ strata $n_i = \frac{100}{4} = 25$ units in each stratum.

Disadvantages of equal allocation:

May need to use weighting to have unbiased estimates

Proportional allocation

Make the proportion of each stratum sampled identical to the proportion of the population. Ie Let the sample fraction be $f = n/N$. So, $n_i = fN_i = n \frac{N_i}{N}$, Where $\frac{N_i}{N}$ is the stratum weight.

Example: $N = 1000$, $n = 100$ $f = \frac{100}{1000} = 0.1$ now suppose $N_1 = 700$ and $N_2 = 300$ then $n_1 = 700 * 0.1 = 70$ and $n_2 = 300 * 0.1 = 30$

Disadvantage of proportional allocation:

Sample size in a stratum may be low thus providing unreliable stratum-specific results.

d)..Cluster Sampling

In many practical situations the population elements are grouped into a number of clusters. A list of clusters can be constructed as the sampling frame but a complete list of elements is often unavailable, or too expensive to construct. In this case it is necessary to use cluster sampling where a random sample of clusters is taken and some or all elements in the selected clusters are observed. Cluster sampling is also preferable in terms of cost, because it is much cheaper, easier and quicker to collect data from adjoining elements than elements chosen at random. On the other hand, cluster sampling is less informative and less efficient per elements in the sample, due to similarities of elements within the same cluster. The loss of efficiency, however, can often be compensated by increasing the overall sample size. Thus, in terms of unit cost, the cluster sampling plan is efficient.

e)..Multi-Stage Samples

Here the respondents are chosen through a process of defined stages. Eg residents within Kibera (Nairobi) may have been chosen for a survey through the following process:

Throughout the country (Kenya) the Nairobi may have been selected at random, (stage 1), within Nairobi, Langata (constituency) is selected again at random (stage 2), Kibera is then selected within Langata (stage 3), then polling stations from Kibera (stage 4) and then individuals from the electoral voters' register (stage 5)! As demonstrated five stages were gone through before the final selection of respondents were selected from the electoral voters' register.

Advantages of probability sample

- (i) Provides a quantitative measure of the extent of variation due to random effects
- (ii) Provides data of known quality
- (iii) Provides data in timely fashion
- (iv) Provides acceptable data at minimum cost
- (v) Better control over nonsampling sources of errors
- (vi) Mathematical statistics and probability can be applied to analyze and interpret the data

1.4.4 Non-probability Sampling

Social research is often conducted in situations where a researcher cannot select the kinds of probability samples used in large-scale social surveys. For example, say you wanted to study homelessness - there is no list of homeless individuals nor are you likely to create such a list. However, you need to get some kind of a sample of respondents in order to conduct your research. To gather such a sample, you would likely use some form of non-probability sampling.

There are four primary types of non-probability sampling methods:

a)..Convenience Sampling

It's a method of choosing subjects who are available or easy to find. This method is also sometimes referred to as haphazard, accidental, or availability sampling. The primary advantage of the method is that it is very easy to carry out, relative to other methods.

Demerit

- One can never be certain what population the participants in the study represent. The population is unknown.
- The method is haphazard, and the cases studied probably don't represent any population you could come up with. However, it's very useful for pilot studies

Advantages of convenience sample

- (i) It's *very easy* to carry out with few rules governing how the *sample* should be collected.
- (ii) The *relative cost* and *time* required to carry out a convenience sample are *small* in comparison to probability sampling techniques. This enables you to achieve the *sample size* you want in a *relatively fast* and *inexpensive* way.
- (iii) The convenience sample may help you gather useful data and information that would not have been possible using *probability sampling techniques*, which require more formal access to *lists of populations* [see, for example, the article on simple random sampling].

For example, imagine you were interested in understanding more about employee satisfaction in a single, large organisation in the United States. You intended to collect your data using a questionnaire. The manager who has kindly given you access to conduct your research is unable to get permission to get a *list* of all employees in the organisation, which you would need to use a *probability sampling technique* such as simple random sampling or systematic random sampling.

However, the manager has managed to secure permission for you to spend two days in the organisation to collect as many questionnaire responses as possible. You decide to spend the two days at the entrance of the organisation where all employees have to pass through to get to their desks. Whilst a *probability sampling technique* would have been preferred, the convenience sample was the only sampling technique that you could use to collect data. Irrespective of the disadvantages of convenience sampling, discussed below, without the use of this sampling technique, you may not have been able to get access to any data on employee satisfaction in the organisation.

Disadvantages of convenience sampling

- The convenience sample often suffers from a number of *biases*. This can be seen in both of our examples, whether the 10,000 students we were studying, or the employees at the large organisation. In both cases, a convenience sample can lead to the *under-representation* or *over-representation* of particular *groups* within the *sample*. If we take the large organisation:

It may be that the organisation has multiple sites, with employee satisfaction varying considerably between these sites. By conducting the survey at the headquarters of the organisation, we may have missed the differences in employee satisfaction amongst those at different sites, including non-office workers. We also do not know why some employees agreed to take part in the survey, whilst others did not. Was it because some employees were simply too busy? Did they not trust the intentions of the survey? Did others take part out of kindness or because they had a particular grievance with the organisation? These types of *biases* are quite typical in convenience sampling.

- Since the *sampling frame* is *not known*, and the *sample* is *not chosen at random*, the *inherent bias* in convenience sampling means that the sample is *unlikely* to be *representative* of the *population* being studied. This undermines your ability to make *generalisations* from your *sample* to the *population* you are studying.

If you are an undergraduate or master's level dissertation student considering using *convenience sampling*, you may also want to read more about how to put together your *sampling strategy* [see the section: Sampling Strategy]

b)..Quota Sampling

Quota sampling is designed to overcome the most obvious flaw of availability sampling. Rather than taking just anyone, you set quotas to ensure that the sample you get represents certain characteristics in proportion to their prevalence in the population. Note that for this method, you have to know something about the characteristics of the population ahead of time. Say you want to make sure you have a sample proportional to the population in terms of gender - you have to know what percentage of the population is male and female, then collect sample until yours matches. Marketing studies are particularly fond of this form of research design.

The primary problem with this form of sampling is that even when we know that a quota sample is representative of the particular characteristics for which quotas have been set, we have no way of knowing if sample is representative in terms of any other characteristics. If we set quotas for gender and age, we are likely to attain a sample with good representativeness on age and gender, but one that may not be very representative in terms of income and education or other factors.

Moreover, because researchers can set quotas for only a small fraction of the characteristics relevant to a study quota sampling is really not much better than availability sampling. To reiterate, you must know the characteristics of the entire population to set quotas; otherwise there's not much point to setting up quotas. Finally, interviewers often introduce bias when allowed to self-select respondents, which is usually the case in this form of research. In choosing males 18-25, interviewers are more likely to

choose those that are better-dressed, seem more approachable or less threatening. That may be understandable from a practical point of view, but it introduces bias into research findings.

Imagine that a researcher wants to understand more about the career goals of students at a single university. Let's say that the university has roughly 10,000 students. Suppose we were interested in *comparing the differences* in career goals between *male* and *female* students at the single university. If this was the case, we would want to ensure that the *sample* we selected had a *proportional* number of *male* and *female* students relative to the *population*.

To create a quota sample, there are three steps:

Choose the relevant grouping characteristic and divide the population accordingly *gender*

Calculate a quota (number of *units* that should be included in *each* for group)

Continue to invite units until the quota for each group is met

Advantages of quota sampling

- i) It is particularly useful when you are unable to obtain a probability sample, but you are still trying to create a sample that is as representative as possible of the population being studied. In this respect, it is the non-probability based equivalent of the stratified random sample.
- ii) Unlike probability sampling techniques, especially stratified random sampling, quota sampling is much *quicker* and *easier* to carry out because it does not require a *sampling frame* and the strict use of random sampling techniques.
- iii) The quota sample improves the *representation* of particular *strata* (*groups*) within the *population*, as well as ensuring that these *strata* are *not over-represented*. For example, it would ensure that we have sufficient male students taking part in the research (60% of our *sample size* of 100; hence, 60 male students). It would also make sure we did not have more than 60 male students, which would result in an *over-representation* of male students in our research.
- iv) It allows *comparison* of *groups*.

Disadvantages of quota sampling

- i) In quota sampling, the *sample* has not been chosen using *random selection*, which makes it impossible to determine the possible *sampling error*.
- ii) this *sampling bias*. Thus *nonstatistical inferences* from the *sample* to the *population*. This can lead to problems of *external validity*.
- iii) Also, with quota sampling it must be possible to clearly divide the *population* into *strata*; that is, *each unit* from the population must only belong to *one stratum*. In our example, this would be fairly simple, since our *strata* are *male* and *female* students. Clearly, a student could only be classified as either male or female. No student could fit into both categories (ignoring transgender issues).

c). Purposive Sampling

Purposive sampling is a sampling method in which elements are chosen based on purpose of the study. Purposive sampling may involve studying the entire population of some limited group or a subset of a population. As with other non-probability sampling methods, purposive sampling does not produce a sample that is representative of a larger population, but it can be exactly what is needed in some cases - study of organization, community, or some other clearly defined and relatively limited group.

Advantages of purposive sampling

- i) There are a wide range of *qualitative research designs* that researchers can draw on. Achieving the goals of such qualitative research designs requires different types of *sampling strategy* and *sampling technique*. One of the major benefits of purposive sampling is the wide range of sampling techniques that can be used across such qualitative research designs; purposive sampling techniques that range from *homogeneous sampling* through to *critical case sampling*, *expert sampling*, and more.

- ii) Whilst the various purposive sampling techniques each have different goals, they can provide researchers with the justification to make *generalisations* from the sample that is being studied, whether such generalisations are *theoretical*, *analytic* and/or *logical* in nature. However, since each of these types of purposive sampling differs in terms of the nature and ability to make generalisations, you should read the articles on each of these purposive sampling techniques to understand their relative advantages.
- iii) Qualitative research designs can involve multiple phases, with each phase building on the previous one. In such instances, different types of sampling technique may be required at each phase. Purposive sampling is useful in these instances because it provides a wide range of non-probability sampling techniques for the researcher to draw on. For example, *critical case sampling* may be used to investigate whether a phenomenon is worth investigating further, before adopting an *expert sampling* approach to examine specific issues further.

Disadvantages of purposive sampling

- i) Purposive samples, irrespective of the type of purposive sampling used, *can be* highly prone to *researcher bias*. The idea that a purposive sample has been created based on the *judgement* of the researcher is not a good defence when it comes to alleviating possible researcher biases,
- ii) specially when compared with *probability sampling* techniques that are designed to reduce such biases. However, this judgemental, subjective component of purpose sampling is only a major disadvantage when such judgements are *ill-conceived* or *poorly considered*; that is, where judgements have not been based on clear criteria, whether a theoretical framework, expert elicitation, or some other accepted criteria.
- iii) The subjectivity and non-probability based nature of *unit* selection (i.e. selecting people, cases/organisations, etc.) in purposive sampling means that it can be difficult to defend the representativeness of the sample. In other words, it can be difficult to convince the reader that the judgement you used to select units to study was appropriate. For this reason, it can also be difficult to convince the reader that research using purposive sampling achieved *theoretical/analytic/logical generalisation*. After all, if different units had been selected, would the results and any generalisations have been the same?

d)..Snowball Sampling

Snowball sampling is a method in which a researcher identifies one member of some population of interest, speaks to him/her, and then asks that person to identify others in the population that the researcher might speak to. This person is then asked to refer the researcher to yet another person, and so on.

Snowball sampling is very good for cases where members of a special population are difficult to locate. For example, *populations* that are subject to social stigma and marginalisation, such as suffers of AIDS/HIV, as well as individuals engaged in illicit or illegal activities, including prostitution and drug use. Snowball sampling is useful in such scenarios because:

The method creates a sample with questionable representativeness. A researcher is not sure who is in the sample. In effect snowball sampling often leads the researcher into a realm he/she knows little about. It can be difficult to determine how a sample compares to a larger population. Also, there's an issue of who respondents refer you to - friends refer to friends, less likely to refer to ones they don't like, fear, etc.

Snowball sampling is a useful choice of *sampling strategy* when the *population* you are interested in studying is *hidden* or *hard-to-reach*.

Advantages of Snowball Sampling

- (i) The chain referral process allows the researcher to reach populations that are difficult to sample when using other sampling methods.
- (ii) The process is cheap, simple and cost-efficient.
- (iii) This sampling technique needs little planning and fewer workforce compared to other sampling techniques.

Disadvantages of Snowball Sampling

- (i) The researcher has little control over the sampling method. The subjects that the researcher can obtain rely mainly on the previous subjects that were observed.
- (ii) Representativeness of the sample is not guaranteed. The researcher has no idea of the true distribution of the population and of the sample.
- (iii) Sampling bias is also a fear of researchers when using this sampling technique. Initial subjects tend to nominate people that they know well. Because of this, it is highly possible that the subjects share the same traits and characteristics, thus, it is possible that the sample that the researcher will obtain is only a small subgroup of the entire population

1.4.5 Limitations of Sampling

- a) Sampling frame: may need complete enumeration
- b) Errors of sampling may be high in small areas
- c) May not be appropriate for the study objectives/questions
- d) Representativeness may be vague, controversial

1.4.6 Characteristics of Good sampling

A good sample should;

- a) Meet the requirements of the study objectives
- b) Provides reliable results
- c) Clearly understandable
- d) Manageable/realistic: could be implemented
- e) Time consideration: reasonable and timely
- f) Cost consideration: economical
- g) Interpretation: accurate, representative
- h) Acceptability

1.5 Survey Administration

1.5.1 Steps in Survey

1. Setting the study objectives; What are the objectives of the study? Is survey the best procedure to collect data? Why other study design (experimental, quasi-experimental, community randomized trials, epidemiologic designs, e.g., case-control study) is not appropriate for the study? What information/data need to be collected?

2. Defining the study *population*; Representativeness Sampling frame

3. Decide sample design: alternative considerations

4. Questionnaire design; Appropriateness, acceptability, culturally appropriate, understandable Pre-test: Appropriate, acceptable, culturally appropriate, will answer

5. Fieldwork; Training/Supervision Quality monitoring Timing: seasonality

6. Quality assurance Every steps Minimizing errors/bias/cheating

7. Data entry/compilation Validation Feedback

8. Analysis: Design consideration

9. Dissemination

10. Plans for next survey: what did you learn, what did you miss?

1.5.2 Modes of Survey Administration

- a) Self-Administered Surveys
- b) Personal interview
- c) Telephone
- d) Mail
- e) Computer assisted self-interviewing(CASI) Variants: CAPI (personal interview); CATI (telephone interview) – Replaces the papers
- f) Combination of methods

a)..Self-Administered Surveys

Self-administered surveys have special strengths and weaknesses.

They are useful in describing the characteristics of a large population and make large samples feasible.

Advantages:

- i) **Low cost.** Extensive training is not required to administer the survey. Processing and analysis are usually simpler and cheaper than for other methods.
- ii) **Reduction in biasing error.** The questionnaire reduces the bias that might result from personal characteristics of interviewers and/or their interviewing skills.
- iii) **Greater anonymity.** Absence of an interviewer provides greater anonymity for the respondent. This is especially helpful when the survey deals with sensitive issues.
- iv) Convenience to the respondents (may complete any time at his/her own convenient time)
- v) Accessibility (greater coverage, even in the remote areas)
- vi) May provide more reliable information (e.g., may consult with others or check records to avoid recall bias)

Disadvantages:

- i) **Requires simple questions.** The questions must be straightforward enough to be comprehended solely on the basis of printed instructions and definitions.
- ii) **No opportunity for probing.** The answers must be accepted as final. Researchers have no opportunity to clarify ambiguous answers.
- iii) **Low response rate;** respondents may not respond to all questions and/or may not return questionnaire
- iv) The respondent must be literate to read and understand the questionnaire
- v) Introduce self selection bias
- vi) Not suitable for complex questionnaire

b). . Interview Surveys

Unlike questionnaires interviewers ask questions orally and record respondents' answers. This type of survey generally decreases the number of —"do not know" and —"no answer" responses, compared with self-administered surveys. Interviewers also provide a guard against confusing items. If a respondent has misunderstood a question, the interviewer can clarify, thereby obtaining relevant responses.

Interviewer selection: background characteristics (race, sex, education, culture) listening skill recording skill experience unbiased observation/recording

Interviewer training: be familiar with the study objectives and significance thorough familiarity with the questionnaire contextual and cultural issues privacy and confidentiality informed consent and ethical issues unbiased view mock interview session

Supervision of the interviewer: Spot check Questionnaire check Reinterview (reliability check)

Advantages

- i) **Flexibility.** Allows flexibility in the questioning process and allows the interviewer to clarify terms that are unclear.
- ii) **Control of the interview situation.** Can ensure that the interview is conducted in private, and respondents do not have the opportunity to consult one another before giving their answers.
- iii) **High response rate.** Respondents who would not normally respond to a mail questionnaire will often respond to a request for a personal interview.
- iv) May record non-verbal behaviour, activities, facilities, contexts
- v) Complex questionnaire may be used
- vi) Illiterate respondents may participate

Disadvantages

- i) ***Higher cost.*** Costs are involved in selecting, training, and supervising interviewers; perhaps in paying them; and in the travel and time required to conduct interviews.
- ii) ***Interviewer bias.*** The advantage of flexibility leaves room for the interviewer's personal influence and bias, making an interview subject to interviewer bias.
- iii) ***Lack of anonymity.*** Often the interviewer knows all or many of the respondents. Respondents may feel threatened or intimidated by the interviewer, especially if a respondent is sensitive to the topic or to some of the questions.
- iv) Less accessibility
- v) Inconvenience
- vi) Often no opportunity to consult records, families, relatives

c).. Telephone Interview

Advantages:

- (i) Less expensive
- (ii) Shorter data collection period than personal interviews
- (iii) Better response than mail surveys

Disadvantages:

- (i) Biased against households without telephone, unlisted number
- (ii) Nonresponse
- (iii) Difficult for sensitive issues or complex topics
- (iv) Limited to verbal responses

d)... Focus Groups

Focus groups are useful in obtaining a particular kind of information that would be difficult to obtain using other methodologies. A focus group typically can be defined as a group of people who possess certain characteristics and provide information of a qualitative nature in a focused discussion

Focus groups generally are composed of six to twelve people. Size is conditioned by two factors: the group must be small enough for everyone to participate, yet large enough to provide diversity. This group is special in terms of purpose, size, composition, and procedures. Participants are selected because they have certain characteristics in common that relate to the topic at hand, such as parents of gang members, and, generally, the participants are unfamiliar with each other. Typically, more than one focus group should be convened, since a group of seven to twelve people could be too atypical to offer any general insights on the gang problem.

A trained moderator probes for different perceptions and points of view, without pressure to reach consensus. Focus groups have been found helpful in assessing needs, developing plans, testing new ideas, or improving existing programs

Advantages:

- i) Flexibility allows the moderator to probe for more in-depth analysis and ask participants to elaborate on their responses.
- ii) Outcomes are quickly known.
- iii) They may cost less in terms of planning and conducting than large surveys and personal interviews.

Limitations

- i) A skilled moderator is essential.
- ii) Differences between groups can be troublesome to analyze because of the qualitative nature of the data.
- iii) Groups are difficult to assemble. People must take the time to come to a designated place at a particular time.

- iv) Participants may be less candid in their responses in front of peers.

New Terminology in Computer Age

PPI: Paper and Pencil Interview

CAI: Computer-Assisted Interview

CATI: Computer-Assisted Telephone Interview

CAPI: Computer-Assisted Personal Interview

CASI: Computer-Assisted Self-Interview

Internet Surveys: Surveys over the WWW

1.6 Sample Size Determination

Sample Size Determination is influenced factors like the purpose of the study, population size, the risk of selecting a "bad" sample, and the allowable sampling error.

There are several approaches to determining the sample size. These include using a census for small populations, imitating a sample size of similar studies, using published tables, and applying formulas to calculate a sample size.

Using a census for small populations

One approach is to use the entire population as the sample. It's impractical for large populations. A census eliminates sampling error and provides data on all the individuals in the population. Finally, virtually the entire population would have to be sampled in small populations to achieve a desirable level of precision

Using a sample size of a similar study

Another approach is to use the same sample size as those of studies similar to the one you plan. Without reviewing the procedures employed in these studies you may run the risk of repeating errors that were made in determining the sample size for another study. However, a review of the literature in your discipline can provide guidance about "typical" sample sizes which are used.

Using published tables

One can also rely on published tables which provide the sample size for a given set of criteria. Yamane, 1967 Table 2.1 and Table 2.2 present sample sizes that would be necessary for given combinations of precision, confidence levels, and variability.

NB, i) these sample sizes reflect the number of *obtained* responses, and not necessarily the number of surveys mailed or interviews planned (this number is often increased to compensate for non-response).
ii) the sample sizes in Table 2.2 presume that the attributes being measured are distributed normally or nearly so. If this assumption cannot be met, then the entire population may need to be surveyed.

Using formulas to calculate a sample size

Yamane (1967) provides a simplified formula to calculate sample sizes. A 95% confidence level and $P = .5$ are assumed for this Equation. $n = \frac{N}{1 + Ne^2}$ Where n is the sample size, N is the population size, and e is the level of precision

2. Data Presentation

2.1 Frequency Distributions Tables

Definitions

Raw Data: unprocessed data ie data in its original form.

Frequency Distribution: The organization of raw data in table form with classes and frequencies. Rather it's a list of values and the number of times they appear in the data set. We have grouped and ungrouped frequency distribution tables for large and small data sets respectively.

2.1.1 Construction of Ungrouped Frequency Distributions

- Note the largest and smallest observations in the data
- Starting with the smallest value, tally the observations of each quantity.
- Count the number of tallies for each quantity and record it as frequency.

Example: Construct an ungrouped frequency table for the data below

16 14 15 13 12 14 16 15 15 14 17 16 13 16 15 14 18 13 15 17

Solution

Value	12	13	14	15	16	17	18
Tally	/	///	////	////	////	//	/
Frequency	1	3	4	5	4	2	1

Note ; when tallying /// is used for 5 counts and not ////

2.1.2 Construction of Grouped Frequency Distributions

When the number of observations is too large and/or when the variable of interest is continuous, it's cumbersome to consider the repetition of each observation. A quick and more convenient way is to group the range of values into a number of exclusive groups or classes and count the class frequency. The resulting table is called a grouped frequency distribution table.

A grouped frequency distribution consists of *classes* and their corresponding *frequencies*. Each raw data value is placed into a quantitative or qualitative category called a **class**. The **frequency** of a class then is the number of data values contained in a specific class.

Steps in construction

- Select the number of classes k. Choose the smallest integer k such that

$$2^k > n \Rightarrow k > \frac{\log n}{\log 2}$$
Eg if $n = 30$ $k > \frac{\log 30}{\log 2} \Rightarrow k = 5$
- Identify the largest and smallest observation and compute the range $R = \text{largest} - \text{smallest}$ value.
- Identify the smallest unit of measurement (u) used in the data collection (ie the accuracy of the measurement.)

Eg for the data 10 30 20 50 30 60 $u=10$ For the data 12 15 11 17 13 $u=1$

For the data 1.6 3.2 2.8 5.6 3.5 1.6 $u=0.1$

- Estimate the class width/interval. $i = \text{round up} \left(\frac{R}{k} \right)$ to the nearest u Eg round up 3.13 to the nearest 0.1 is 3.2
- Pick the starting value (lower class limit of the 1st class (LCL_1)) as the smallest value used in the computation of R above. Successive LCL_s are got by adding I to the previous LCL
- Find the upper class limit of the 1st class (UCL_1) by subtracting u from LCL_2 . Successive UCL_s are got by adding I to the previous UCL
- If necessary find the class boundaries as follows $LCB = LCL - \frac{1}{2}u$ and $UCB = UCL + \frac{1}{2}u$
- Tally the number of observations falling in each class and record the frequency.
NB a value x fall into a class $LCL - UCL$ if $LCB \leq x \leq UCB$

Example 1

Organize the data below into a grouped frequency table.

15.0 17.4 10.3 9.2 20.7 18.9 16.6 22.4 23.7 18.6 26.1 16.5 19.7 12.9 15.7
30.8 15.4 20.3 24.0 29.6 18.3 23.7 17.8 24.6 23.0 21.4 32.8 12.5 17.5 18.3
23.2 21.6 20.8 29.8 24.5 28.4 13.5 17.1 27.1 27.9

Solution

$$u = 0.1 \quad n = 40 \Rightarrow k > \frac{\log 40}{\log 2} \Rightarrow k = 6 \quad \text{Range} = 32.8 - 9.2 = 23.6$$

$$i = \text{round up } \left(\frac{23.6}{6}\right) \text{ to the nearest } 0.1 = 4.0$$

$$\text{Now } LCL_1 = 9.2 \Rightarrow LCL_2 = 13.2 \quad LCL_3 = 17.2 \text{ etc}$$

$$\Rightarrow UCL_1 = LCL_2 - u = 13.2 - 0.1 = 13.1, \quad UCL_2 = 17.1 \quad UCL_3 = 21.1 \text{ etc}$$

The frequency table is as shown below

Class	Boundaries	tally	Freq	C.F
9.2 - 13-1	9.15 - 13-15	////	4	4
13.2 - 17.1	13.15 - 17.15	//// //	7	11
17.2 - 21-1	17.15 - 21-15	//// //// /	11	22
21.2 - 25.1	21.15 - 25.15	//// ////	10	32
25.2 - 29.1	25.15 - 29.15	////	4	36
29.2 - 33-	29.15 - 33-15	////	4	40

Sometimes due to convenience, it may be necessary to slightly lower the starting value. Eg in the above case we may use 9.0 in place of 9.2.

Example 2

These data represent the record high temperatures in degrees Fahrenheit (F) for each of the 50 states. Construct a grouped frequency distribution for the data.

112 100 127 120 134 118 105 110 109 112 110 118 117 116 118 122 114 114 105
109 107 112 114 115 118 117 118 122 106 110 116 108 110 121 113 120 119 111
104 111 120 113 120 117 105 110 118 112 114 114

Solution

$$u = 1 \quad n = 50 \Rightarrow k > \frac{\log 50}{\log 2} \Rightarrow k = 6 \quad \text{Range} = 134 - 100 = 34$$

$$i = \text{round up } \left(\frac{34}{6}\right) \text{ to the nearest } 1 = 6$$

$$\text{Now } LCL_1 = 100 \Rightarrow LCL_2 = 106 \quad LCL_3 = 112 \text{ etc}$$

$$\Rightarrow UCL_1 = LCL_2 - u = 106 - 1 = 105, \quad UCL_2 = 111 \quad UCL_3 = 117 \text{ etc}$$

The frequency table is as shown below

Class	Boundaries	tally	Freq	C.F
100 - 105	99.5 - 105.5	////	5	5
106 - 111	105.5 - 111.5	//// //// ///	13	18
112 - 117	111.5 - 117.5	//// //// //// //	16	34
118 - 123	117.5 - 123.5	//// //// ////	14	48
124 - 129	124.5 - 129.5	/	1	49
130 - 135	129.5 - 135-5	/	1	50

Exercise

- 1) The data shown here represent the number of miles per gallon (mpg) that 30 selected four-wheel-drive sports utility vehicles obtained in city driving. Construct a frequency distribution, and analyze the distribution. 12 17 12 14 16 18 16 18 12 16 17 15 15 16 12 15 16 16 12 14 15 12 15 15 19 13 16 18 16 14

- 2) Suppose a researcher wished to do a study on the ages of the top 50 wealthiest people in the world. The researcher first would have to get the data on the ages of the people. In this case, these ages are listed in *Forbes Magazine*. 49 57 38 73 81 74 59 76 65 69 54 56 69 68 78 65 85 49 69 61 48 81 68 37 43 78 82 43 64 67 52 56 81 77 79 85 40 85 59 80 60 71 57 61 69 61 83 90 87 74 Organize the data into a grouped frequency table
- 3) The data represent the ages of our Presidents at the time they were first inaugurated. 57 61 57 57 58 57 61 54 68 51 49 64 50 48 65 52 56 46 54 49 50 47 55 55 54 42 51 56 55 54 51 60 62 43 55 56 61 52 69 64 46 54
 - a) Were the data obtained from a population or a sample? Explain your answer.
 - b) What was the age of the oldest and youngest President?
 - c) Construct a frequency distribution for the data.
 - d) Are there any peaks in the distribution?
 - e) identify any possible outliers.
- 4) The state gas tax in cents per gallon for 25 states is given below. Construct a grouped frequency distribution for the data. 7.5 16 23.5 17 22 21.5 19 20 27.1 20 22 20.7 17 28 20 23 18.5 25.3 24 31 14.5 25.9 18 30 31.5
- 5) Listed are the weights of the NBA's top 50 players. Construct a grouped frequency distribution and analyze the results in terms of peaks, extreme values, etc. 240 210 220 260 250 195 230 270 325 225 165 295 205 230 250 210 220 210 230 202 250 265 230 210 240 245 225 180 175 215 215 235 245 250 215 210 195 240 240 225 260 210 190 260 230 190 210 230 185 260
- 6) The number of stories in each of the world's 30 tallest buildings is listed below. Construct a grouped frequency distribution and analyze the results in terms of peaks, extreme values, etc. 88 88 110 88 80 69 102 78 70 55 79 85 80 100 60 90 77 55 75 55 54 60 75 64 105 56 71 70 65 72
- 7) The average quantitative GRE scores for the top 30 graduate schools of engineering are listed. Construct a grouped frequency distribution and analyze the results in terms of peaks, extreme values, etc. 767 770 761 760 771 768 776 771 756 770 763 760 747 766 754 771 771 778 766 762 780 750 746 764 769 759 757 753 758 746
- 8) The number of passengers (in thousands) for the leading U.S. passenger airlines in 2004 is indicated below. Use the data to construct a grouped frequency distribution with a reasonable number of classes and comment on the shape of the distribution.
- 9) 91.570 86.755 81.066 70.786 55.373 42.400 40.551 21.119 16.280 14.869 13.659 13.417 3.170 12.632 11.731 10.420 10.024 9.122 7.041 6.954 6.406 6.362 5.930 5.585
- 10) The heights (in feet above sea level) of the major active volcanoes in Alaska are given here. Construct a grouped frequency distribution for the data. 4,265 3,545 4,025 7,050 11,413 3,490 5,370 4,885 5,030 6,830 4,450 5,775 3,945 7,545 8,450 3,995 10,140 6,050 10,265 6,965 150 8,185 7,295 2,015 5,055 5,315 2,945 6,720 3,465 1,980 2,560 4,450 2,759 9,430 7,985 7,540 3,540 11,070 5,710 885 8,960 7,015

2.2 Graphical Displays

After you have organized the data into a frequency distribution, you can present them in graphical form. The purpose of graphs in statistics is to convey the data to the viewers in pictorial form. It is easier for most people to comprehend the meaning of data presented graphically than data presented numerically in tables or frequency distributions. This is especially true if the users have little or no statistical knowledge.

Statistical graphs can be used to describe the data set or to analyze it. Graphs are also useful in getting the audience's attention in a publication or a speaking presentation. They can be used to discuss an issue, reinforce a critical point, or summarize a data set. They can also be used to discover a trend or pattern in a situation over a period of time.

The commonly used graphs in research are; the pie chart, bar chart, histogram, frequency polygon and the cumulative frequency curve (Ogive).

2.2.1 Pie Chart

It's a circular graph having radii divide a circle into sectors proportional in angle to the relative size of the quantities in the category being represented.

How to Draw

- (i) Add up the given quantities and let s be the sum of the values
- (ii) For each quantity x , calculate the representative angle and percentage as $\frac{x}{s}(360^\circ)$ and $\frac{x}{s}(100\%)$ respectively
- (iii) Draw a circle and divide it into sectors using the angles calculated in step ii above
- (iv) Label the sector by the group represented and indicate the corresponding percentage.

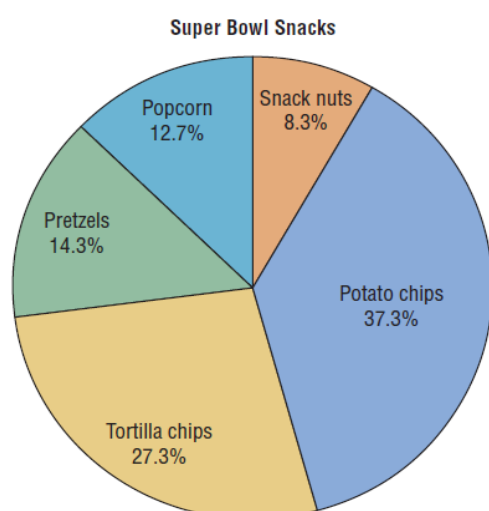
Example

This frequency distribution shows the number of pounds of each snack food eaten during the Super Bowl. Construct a pie graph for the data.

Snack	Potato chips	Tortilla chips	Pretzels	Popcorn	Snack nuts
Pounds (in millions)	11.2	8.2	4.3	3.8	2.5

Solution

Snack	Potato chips	Tortilla chips	Pretzels	Popcorn	Snack nuts	Total
Pounds (in millions)	11.2	8.2	4.3	3.8	2.5	30.0
Representative Angle	134	98	52	46	30	360
Representative %age	37.3	27.3	14.3	12.7	8.3	99.9



2.2.2 Bar chart

A bar chart consists of a set of equal spaced rectangles whose heights are proportional to the frequency of the category /item being considered. The X axis in a bar chart can represent the number of categories.

Note: Bars are of uniform width and there is equal spacing between the bars.

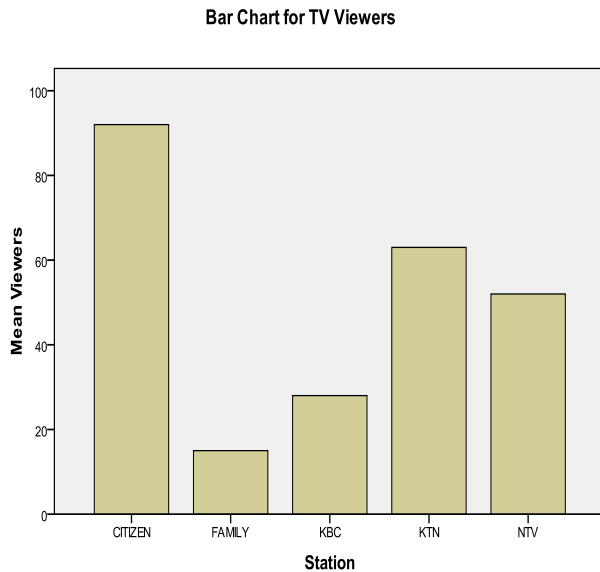
Example

A sample of 250 students was asked to indicate their favourite TV channels and their responses were as follows.

TV station	KBC	NTV	CITIZEN	KTN	FAMILY
No. of viewers	28	52	92	63	15

Draw a bar chart to represent this information.

Solution



2.23 Pareto Charts

It consists of a set of continuous rectangles where the variable displayed on the horizontal axis is qualitative or categorical and the frequencies are displayed by the heights of vertical bars, which are arranged in order from highest to lowest. A **Pareto chart** is used to represent a frequency distribution for a categorical variable.

Points to note when drawing a Pareto Chart

- Make the bars the same width.
- Arrange the data from largest to smallest according to frequency.
- Make the units that are used for the frequency equal in size.

When you analyze a Pareto chart, make comparisons by looking at the heights of the bars.

Example

The table shown here is the average cost per mile for passenger vehicles on state turnpikes. Construct and analyze a Pareto chart for the data.

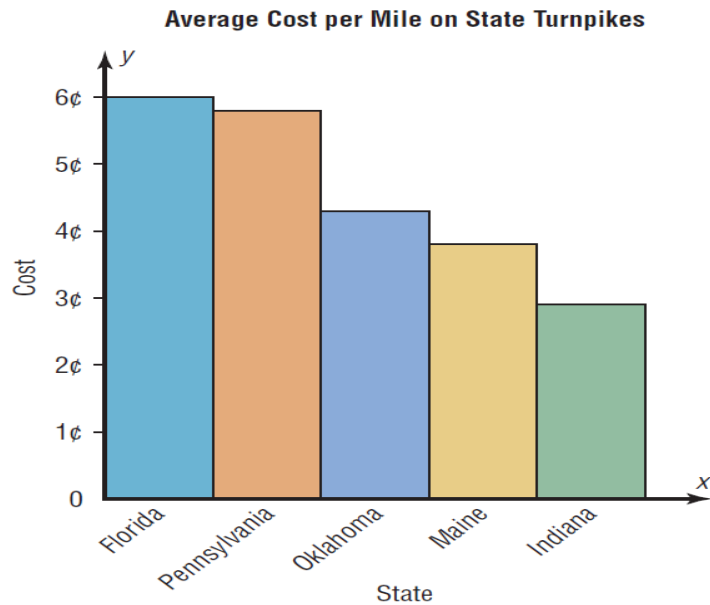
State	Indiana	Oklahoma	Florida	Maine	Pennsylvania
Number	2.9	4.3	6.0	3.8	5.8

Solution

Arrange the data from the largest to smallest according to frequency.

State	Florida	Pennsylvania	Oklahoma	Maine	Indiana
Number	6.0	5.8	4.3	3.8	2.9

Draw and label the x and y axes and then the bars corresponding to the frequencies. The Pareto chart shows that Florida has the highest cost per mile. The cost is more than twice as high as the cost for Indiana.



2.2.4 Histogram

It consists of a set of continuous rectangles such that the areas of the rectangles are proportional to the frequency. For ungrouped data, the heights of each bar is proportional to frequency. For grouped data, the height of each rectangle is the relative frequency h and is given by $h = \frac{\text{frequency}}{\text{IClass Interval}}$. The width of the bars is determined by the class boundaries.

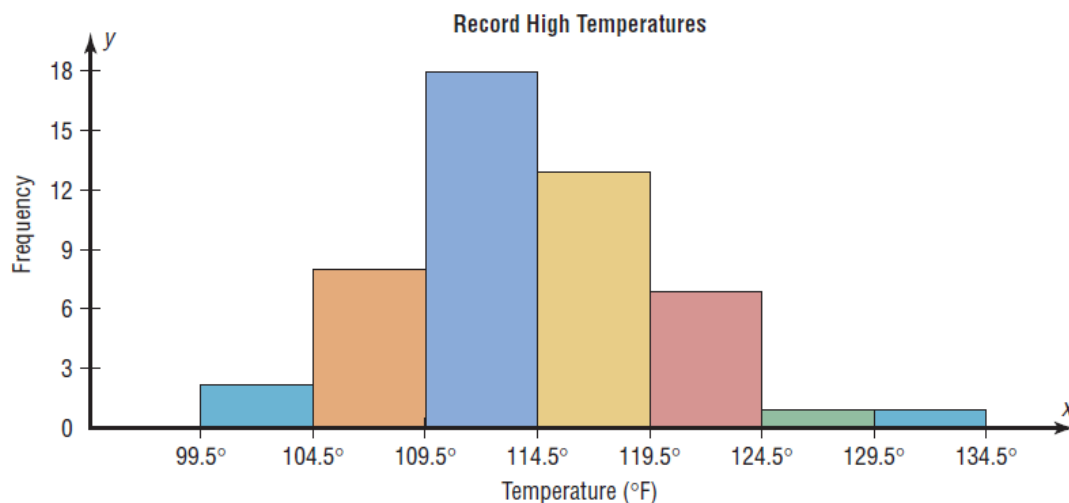
Example

Construct a histogram and an ogive to represent the data shown below

Class	100-104	105-109	110 -114	115-119	120 - 24	125-129	130 -134
Freq	2	8	18	13	7	1	1

Solution

Boundaries	99.5-104.5	104.5-109.5	109.5 - 114.5	114.5-119.5	119.5 – 124.5	124.5-129.5	129.5 - 134.5
Heights	2	8	18	13	7	1	1
CF	2	10	28	41	48	49	50



2.2.5 Frequency polygon

It's a plot of frequency against mid points joined with straight line segments between consecutive points. It can also be obtained by joining the mid point of the tops of the bars in a histogram. The gaps at both ends are filled by extending to the next lower and upper imaginary classes assuming frequency zero.

Example

Fit a frequency polygon on your histogram above.

Example: Consider the following frequency distribution.

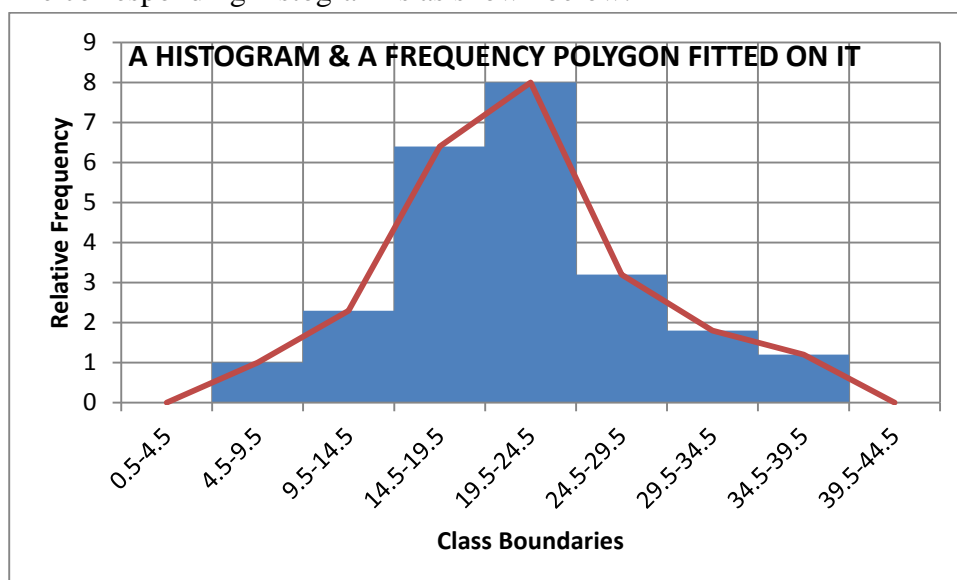
Class	5-9	10-14	15-19	20-24	25-29	30-34	35-39
freq	5	12	32	40	16	9	6

Draw a histogram to represent this information and fit a frequency polygon on it.

Solution

Boundaries	4.5-9.5	9.5-14.5	14.5-19.5	19.5-24.5	24.5-29.5	29.5-34.5	34.5-39.5
heights	1	2.4	6.4	8	3.2	1.8	1.2

The corresponding histogram is as shown below.



2.2.6 Cumulative frequency curve (ogive)

It is a plot of cumulative frequency against upper boundaries joined with a smooth curve. The gap on the lower end is filled by extending to the next lower imaginary class assuming frequency zero. This graph is useful in estimating median and other measures of location.

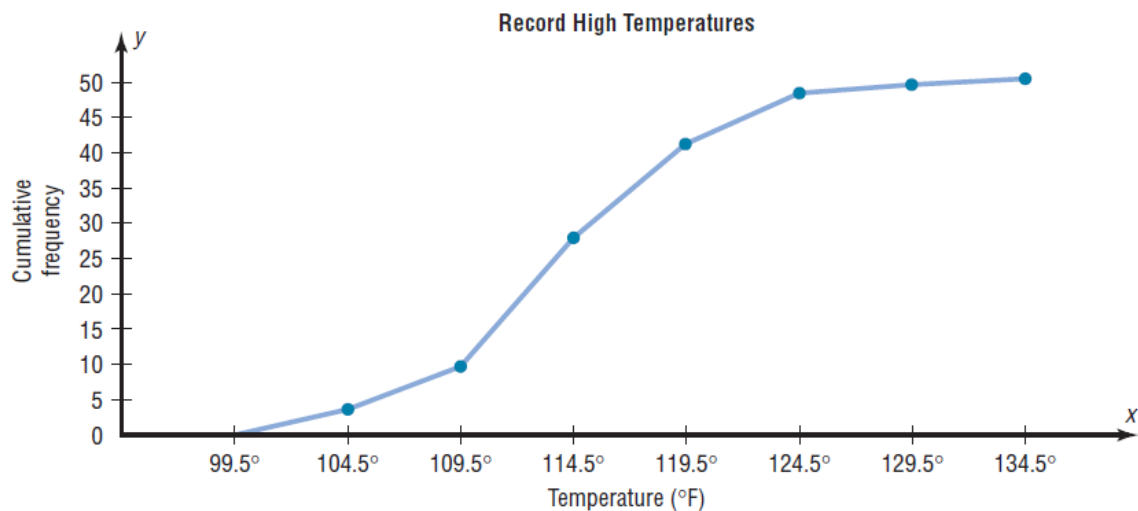
Example:

Construct an ogive to represent the data shown below

Class	100-104	105-109	110 -114	115-119	120 - 24	125-129	130 -134
Freq	2	8	18	13	7	1	1

Solution

Boundaries	99.5-104.5	104.5-109.5	109.5 - 114.5	114.5-119.5	119.5 – 124.5	124.5-129.5	129.5 - 134.5
CF	2	10	28	41	48	49	50



Exercise

- Construct a pie chart and a bar graph showing the blood types of the army inductees described in the frequency distribution is repeated here.

Blood group	A	B	AB	O
Frequency	5	7	4	9

- The table below shows the average money spent by first-year college students. Draw a pie chart and a bar graph for the data.

Nature of Expense	Electronics	Dorm decor	Clothing	Shoes
Amount(in \$)	728	344	141	72

- The table shown here is the average cost per mile for passenger vehicles on state turnpikes. Draw a pie chart and a bar graph for the data.

State	Indiana	Oklahoma	Florida	Maine	Pennsylvania
Number	2.9	4.3	6.0	3.8	5.8

- The following data are based on a survey from American Travel Survey on why people travel. Construct a pie chart a bar graph and a Pareto chart for the data and analyze the results.

Purpose	Personal business	Visit friends or relatives	Work-related	Leisure
Number	146	330	225	299

- The following percentages indicate the source of energy used worldwide. Construct a Pareto chart and a vertical pie chart, a bar graph and a Pareto graph for the energy used.

Energy Type	Petroleum	Coal	Dry natural gas	Hydroelectric	Nuclear	Others
percentage	39.8	23.2	22.4	7.0	6.4	1.2

- The following elements comprise the earth's crust, the outermost solid layer. Illustrate the composition of the earth's crust with a pie chart and a bargraph for this data.

Element	Oxygen	Silicon	Aluminum	Iron	Calcium	Others
percentage	45.6	27.3	8.4	6.2	4.7	7.8

- The sales of recorded music in 2004 by genre are listed below. Represent the data with an appropriate graph.

Rock	Country	Rap/hip-hop	R&B/urban	Pop	Religious	Children's	Jazz	Classical	Oldies	Soundtracks	New age	Others
23.9	13.0	12.1	11.3	10.0	6.0	2.8	2.7	2.0	1.4	1.1	1.0	8.9

- 8) The top 10 airlines with the most aircraft are listed. Represent these data with an appropriate graph.

American	Continental	United	Southwest	Northwest	American Eagle	U.S. Airways	Lufthansa (Ger.)
714	364	603	327	424	245	384	233

- 9) The top prize-winning countries for Nobel Prizes in Physiology or Medicine are listed here. Represent the data with an appropriate graph.

United States	Denmark	United Kingdom	Austria	Germany	Belgium	Sweden	Italy	France	Australia	Switzerland
80	5	24	4	16	4	8	3	7	3	6

- 10) Construct a histogram, frequency polygon, and an ogive for the distribution (shown here) of the miles that 20 randomly selected runners ran during a given week.

Class	6-10	11-15	16 -20	21-25	26 - 30	31-35	36 -40
Freq	1	2	3	5	4	3	2

- 11) For 108 randomly selected college applicants, the following frequency distribution for entrance exam scores was obtained. Construct a histogram, frequency polygon, and ogive for the data.

Class	90-98	99-107	108-116	117-125	126-134
Freq	6	22	43	28	9

Applicants who score above 107 need not enrol in a summer developmental program.

In this group, how many students do not have to enroll in the developmental program?

- 12) Thirty automobiles were tested for fuel efficiency, in miles per gallon (mpg). The following frequency distribution was obtained. Construct a histogram, a frequency polygon, and an ogive for the data.

Class	8-12	13-17	18-22	23-27	28-32
Freq	3	5	15	5	2

- 13) The salaries (in millions of dollars) for 31 NFL teams for a specific season are given in this frequency distribution.

Class	39.9-42.8	42.9-45.8	45.9-48.8	48.9-51.8	51.9-54.8	54.9-57.8
Freq	2	2	5	5	12	5

Construct a histogram, a frequency polygon, and an ogive for the data; and comment on the shape of the distribution.

- 14) In a study of reaction times of dogs to a specific stimulus, an animal trainer obtained the following data, given in seconds. Construct a histogram, a frequency polygon, and an ogive for the data; analyze the results.

Class	2.3-2.9	3.0-3.6	3.7-4.3	4.4-5.0	5.1-5.7	5.8-6.4
Freq	10	12	6	8	4	2

- 15) The animal trainer in question above selected another group of dogs who were much older than the first group and measured their reaction times to the same stimulus. Construct a histogram, a frequency polygon, and an ogive for the data.

Class	2.3-2.9	3.0-3.6	3.7-4.3	4.4-5.0	5.1-5.7	5.8-6.4
Freq	1	3	4	16	14	4

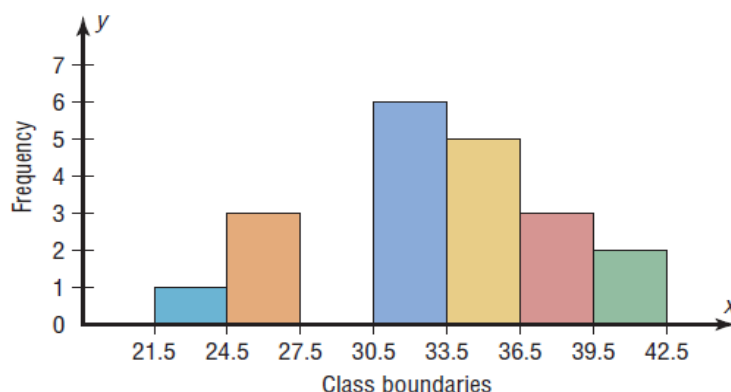
Analyze the results and compare the histogram for this group with the one obtained in the above question. Are there any differences in the histograms?

- 16) The frequency distributions shown indicate the percentages of public school students in fourth-grade reading and mathematics who performed at or above the required proficiency levels for the 50 states in the United States. Draw histograms for each, and decide if there is any difference in the performance of the students in the subjects.

Class	18-22	23-27	28-32	33-37	38-42	43-48
Reading Freq	7	6	14	19	3	1
Math Freq	5	9	11	16	8	1

Using the histogram shown here, Construct a frequency distribution; include class limits, class frequencies, midpoints, and cumulative frequencies. Hence answer these questions.

- How many values are in the class 27.5–30.5?
- How many values fall between 24.5 and 36.5?
- How many values are below 33.5?
- How many values are above 30.5?



2.3 Stem and Leaf Plots

The stem and leaf plot is a method of organizing data and is a combination of sorting and graphing. It has the advantage over a grouped frequency distribution of retaining the actual data while showing them in graphical form. A **stem and leaf plot** is a data plot that uses part of the data value as the stem and part of the data value as the leaf to form groups or classes. For this plot, it's easy to identify the mode, the smallest value and the largest value.

Note:

- In a stem and leaf plot, classes width/interval must be uniform.
- The leaves in the final stem and leaf plot should be arranged in order.

Example 1

At an outpatient testing center, the number of cardiograms performed each day for 20 days is shown. Construct a stem and leaf plot for the data.

25 31 20 32 13 14 43 02 57 23 36 32 33 32 44 32 52 44 51 45

Solution

Arrange the data in order and separate the data according to the first digit, as shown.

02 13, 14 20, 23, 25 31, 32, 32, 32, 32, 33, 36 43, 44, 44, 45 51, 52, 57

Remark: Arranging the data in order is not essential and can be cumbersome when the data set is large; however, it is helpful in constructing a stem and leaf plot.

A display can be made by using the leading digit as the *stem* and the trailing digit as the *leaf*. For example, for the value 32, the leading digit, 3, is the stem and the trailing digit, 2, is the leaf. For the value 14, the 1 is the stem and the 4 is the leaf. Now a plot can be constructed as shown in below.

10's digit (stem)	1's digit (leaf)
0	2
1	3, 4
2	0, 3, 5
3	1, 2, 2, 2, 2, 3, 6
4	3, 4, 4, 5
5	1, 2, 7

The plot above shows that the distribution peaks in the center and that there are no gaps in the data. For 7 of the 20 days, the number of patients receiving cardiograms was between 31 and 36. The plot also shows that the testing center treated from a minimum of 2 patients to a maximum of 57 patients in any one day.

Note: If there are no data values in a class, you should write the stem number and leave the leaf row blank. Do not put a zero in the leaf row.

When you analyze a stem and leaf plot, look for peaks and gaps in the distribution. See if the distribution is symmetric or skewed. Check the variability of the data by looking at the spread.

Example 2

An insurance company researcher conducted a survey on the number of car thefts in a large city for a period of 30 days last summer. The raw data are shown. Construct a stem and leaf plot by using classes 50–54, 55–59, 60–64, 65–69, 70–74, and 75–79.

52 62 51 50 69 58 77 66 53 57 75 56 55 67 73 79 59 68 65 72 57 51 63 69 75 65 53 78 66 55 Obtain a stem and leaf plot this data.

Solution

Stem	Leaf	Sorted plot	
		Stem	Leaf
5	2, 1, 0, 3, 1, 3	5	0, 1, 1, 2, 3, 3
5	8, 7, 6, 5, 9, 7, 5	5	5, 5, 6, 7, 7, 8, 9
6	2, 3	6	2, 3
6	9, 6, 7, 8, 5, 9, 5, 6	6	5, 5, 6, 6, 7, 8, 9, 9
7	3, 2,	7	2, 3
7	7, 5, 9, 5, 8	7	5, 5, 7, 8, 9

The stem does not necessarily represent the tens digit. See example 3 and 4 below

Example 3

When the data values are in the hundreds, such as 325, the stem is 32 and the leaf is 5. For example, the stem and leaf plot for the data values 325, 327, 330, 332, 335, 341, 345, and 347 looks like this.

Stem	Leaf
32	5, 7
33	0, 2, 5
34	5, 7

Example 4

The maximum temperature measured to the nearest degrees centigrade was recorded in a certain town each day on June. The results were as follows. Draw a stem and leaf diagram with classes 12-14, 15-17, . . . for this data 19, 13, 19, 19, 20, 12, 19, 22, 22, 16, 18, 16, 19, 20, 17, 13, 14, 12, 15, 17, 16, 17, 19, 22, 22, 20, 19, 19, 20, 20

Solution

The class 18-20 cannot be represented by a stem of 1 since the 10's digit changes within the class. Therefore the stem we can use are 12, 15, 18, 21. Here the leaf will be the number in excess to the stem Eg 15/2 means 17 and 18/0 means 18

Stem	Leaf	Arranged plot	
		Stem	Leaf
12	0, 0, 2, 1	12	0, 0, 1, 2
15	1, 2, 0, 2, 1, 1, 2	15	0, 1, 1, 1, 2, 2, 2
18	1, 1, 1, 2, 1, 2, 1, 0, 1, 2, 1, 1, 2, 2	18	0, 1, 1, 1, 1, 1, 1, 1, 2, 2, 2, 2, 2
21	2, 1, 1, 1, 1	21	1, 1, 1, 1, 2

Key 15/2 means 17 and 18/0 means 18

2.3.1 Back-to-Back Stem and Leaf Plot

Related distributions can be compared by using a back-to-back stem and leaf plot.

The back-to-back stem and leaf plot uses the same digits for the stems of both distributions, but the digits that are used for the leaves are arranged in order out from the stems on both sides. Example 2–15 shows a back-to-back stem and leaf plot.

Example

The number of stories in two selected samples of tall buildings in Atlanta and Philadelphia is shown. Construct a back-to-back stem and leaf plot, and compare the distributions.

Atlanta

55 70 44 36 40 63 40 44 34 38
60 47 52 32 32 50 53 32 28 31
52 32 34 32 30 26 29

Philadelphia

61 40 38 32 30 58 40 40 25 30 50 38 36
54 40 36 30 30 53 39 36 34 33 39 32

Solution

The final back-to-back stem and leaf plot looks like the below

Atlanta	Stem	Philadelphia
9, 8, 6	2	5
8, 6, 4, 4, 2, 2, 2, 2, 1	3	0, 0, 0, 0, 2, 2, 3, 4, 6, 6, 6, 8, 8, 9, 9
7, 4, 4, 0, 0	4	0, 0, 0, 0
5, 3, 2, 2, 0, 0	5	0, 3, 4, 8
3, 0	6	1
0	7	

The buildings in Atlanta have a large variation in the number of stories per building. Although both distributions are peaked in the 30- to 39-story class, Philadelphia has more buildings in this class. Atlanta has more buildings that have 40 or more stories than Philadelphia does.

Exercise

- The number of visitors to the Railroad Museum during 24 randomly selected hours is shown here. Construct a stem and leaf plot for the data.
67 62 38 73 34 43 72 35 53 55 58 63 47 42 51 62 32 29 47 62 29 38 36 41
- The numbers of public libraries in operation for selected states are listed below. Organize the data with a stem and leaf plot.
102 176 210 142 189 176 108 113 205 209 184 144 108 192 176
- The age at inauguration for each U.S. President is shown. Construct a stem and leaf plot and analyze the data. 57 54 52 55 51 56 61 68 56 55 54 61 57 51 46 54 51 52 57 49 54 42 60 69 58 64 49 51 62 64 57 48 50 56 43 46 61 65 47 55 55 54
- Write down the list of values represented in the following stem and leaf diagram

Stem	Leaf
5	1, 3
10	2, 4, 4
15	0, 0, 1, 4, 4
20	3, 3, 4
25	1, 1

Key: 5/1 means 6

Stem	Leaf
7	0, 2, 3
11	0, 1, 1, 3, 3, 3
15	1, 2, 2, 3
19	0, 2
23	1

Key: 15/3 means 18

- A listing of calories per one ounce of selected salad dressings (not fat-free) is given below. Construct a stem and leaf plot for the data. 100 130 130 130 110 110 120 130 140 100 140 170 160 130 160 120 150 100 145 145 145 115 120 100 120 160 140 120 180 100 160 120 140 150 190 150 180 160
- A special aptitude test is given to job applicants. The data shown here represent the scores of 30 applicants. Construct a stem and leaf plot for the data and summarize the results.

204 210 227 218 254 256 238 242 253 227 251 243 233 251 241 237 247 211
 222 231 218 212 217 227 209 260 230 228 242 200

- 7) For the following data, draw a stem and leaf plot with classes
 a) 5 - 9, 10 - 14, 15 - 19, . . .
 b) 4 - 6, 7 - 9, 10 - 12, . . . ,
 21, 19, 6, 12, 8, 18, 9, 8, 11, 17, 15, 13, 16, 9, 17, 18, 9, 24, 17, 7, 8, 17, 17,
 8, 7, 11, 16, 17, 8, 5, 13, 22, 20, 16, 20, 13
- 8) Draw a stem and leaf diagram for the following data. use classes 3.0-3.9, 4.0-4.9, . . .
 10.0, 9.2, 7.3, 7.0, 6.5, 5.4, 5.3, 10.1, 8.4, 8.5, 7.1, 7.6, 7.9, 6.9, 9.6, 5.5, 7.4, 7.0, 8.2, 5.5,
 7.8, 8.2, 7.5, 6.1, 6.1, 3.9, 6.8, 7.6, 8.1, 8.0
- 9) For the following stem and leaf diagram give the values in bold and the class containing them if the diagram represent;

Stem	Leaf
0	7
0	9
1	0, 1
1	2, 2
1	4, 4, 4 , 5, 5
1	6, 6, 7, 7
1	8, 8, 8 , 8, 9, 9
2	0, 0 , 1, 1
2	2, 3

- 10) Use a back to back stem and leaf to compare the exam marks for French and English in a class of 20 pupils

French	English
75, 69, 58, 58, 46, 44, 52, 50, 53, 78, 51, 61, 61, 45, 31, 44, 53, 66, 47, 57	52, 58, 68, 77, 88, 85, 43, 44, 56, 65, 65, 79, 44, 71, 84, 72, 63, 69, 72, 79

- 11) The growth (in centimetres) of two varieties of plant after 20 days is shown in this table. Construct a back-to-back stem and leaf plot for the data, and compare the distributions.

Variety 1	Variety 2
20 12 39 38 41 43 51 52 59 55 53 59 50 58 35 38 23 32 43 53	18 45 62 59 53 25 13 57 42 55 56 38 41 36 50 62 45 55

- 12) The math and reading achievement scores from the National Assessment of Educational Progress for selected states are listed below. Construct a back-to-back stem and leaf plot with the data and compare the distributions.

Math	Reading
52 66 69 62 61 63 57 59 59 55 55 59 74 72 73 68 76 73	65 76 76 66 67 71 70 70 66 61 61 69 78 76 77 77 77 80

- 13) Draw a back to back stem and leaf to compare the reaction time of boys and girls to a certain stimuli

Maths	Reading
0.14, 0.19, 0.18, 0.09, 0.19, 0.23, 0.16 0.22, 0.15, 0.16, 0.20, 0.16, 0.16, 0.11 0.15, 0.21, 0.23, 0.22, 0.23, 0.18	0.18, 0.24, 0.16, 0.22, 0.19, 0.19, 0.25, 0.22, 0.21, 0.16, 0.22, 0.18, 0.21, 0.22 0.22, 0.25, 0.17, 0.22, 0.19, 0.19

3 NUMERICAL SUMMARIES

A numerical summary for a set of data is referred to as a statistic if the data set is a sample and a parameter if the data set is the entire population.

Numerical summaries are categorized as measures of location and measures of spread. Measures of location can further be classified into measures of central tendency and measures of relative positioning (quantiles).

3.1 Measures of Location

Before discussing the measures of location, it's important to consider summation notation and indexing

Index (subscript) Notation: Let the symbol x_i (read 'x sub t'i') denote any of the n values x_1, x_2, \dots, x_n assumed by a variable X . The letter i in x_i ($i=1, 2, \dots, n$) is called an index or subscript. The letters j, k, p, q or s can also be used.

Summation Notation: $\sum_{i=1}^n x_i = x_1 + x_2 + \dots + x_n$

Example: $\sum_{i=1}^n X_i Y_i = X_1 Y_1 + X_2 Y_2 + \dots + X_n Y_n$ and

$$\sum_{i=1}^n aX_i = aX_1 + aX_2 + \dots + aX_n = a(X_1 + X_2 + \dots + X_n) = a \sum_{i=1}^n X_i$$

3.1.1 Measures of Central Tendency (Averages)

A Measure of Central Tendency of a set of numbers is a value which best represents it. There are three different types of Central Tendencies namely the mean, median and mode. Each has advantages and disadvantages depending on the data and intended purpose.

Arithmetic Mean

The arithmetic mean of a set of values x_1, x_2, \dots, x_n , denoted \bar{x} if the data set is a sample, is found by dividing the sum of the set of numbers with the actual number of values. I.e

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{x_1 + x_2 + \dots + x_n}{n}$$

Example 1 Find the mean of 1, 2, 3, 4, 5, 6, 7, 8, 9 and 10.

Solution

Sum of values: $1 + 2 + 3 + 4 + 5 + 6 + 7 + 8 + 9 + 10 = 55$

Number of values = 10

Mean of values $\bar{x} = \frac{55}{10} = 5.5$

Note: If the numbers x_1, x_2, \dots, x_n occur f_1, f_2, \dots, f_n times respectively, (occur with frequencies f_1, f_2, \dots, f_n), the arithmetic mean is, given by

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n f_i x_i = \frac{f_1 x_1 + f_2 x_2 + \dots + f_n x_n}{f_1 + f_2 + \dots + f_n}$$

where n is the total frequency. This is the formula for the mean of a grouped data.

Example 2 The grades of a student on six examinations were 84, 91, 72, 68, 91 and 72. Find the arithmetic mean.

The arithmetic mean $\bar{x} = \frac{1}{n} \sum_{i=1}^n f_i x_i = \frac{1(84) + 2(91) + 2(72) + 1(68)}{1 + 2 + 2 + 1} = 79.67$

Example 3 If 5, 8, 6 and 2 occur with frequencies 3, 2, 4 and 1 respectively, the arithmetic mean is $\bar{x} = \frac{1}{n} \sum_{i=1}^n f_i x_i = \frac{3(5) + 2(8) + 4(6) + 1(2)}{3 + 2 + 4 + 1} = 5.7$

Exercise

- 1 Find the mean of 9, 3, 4, 2, 1, 5, 8, 4, 7, 3
- 2 A sample of 5 executives received the following amount of bonus last year: sh 14,000, sh 15,000, sh 17,000, sh 16,000 and sh y. Find the value of y if the average bonus for the 5 executives is sh 15,400

Properties of the Arithmetic Mean

(1) The algebraic sum of the deviations of a set of numbers from their arithmetic mean is

zero, that is $\sum_{i=1}^n (x_i - \bar{x}) = 0$.

(2) $\sum_{i=1}^n (x_i - a)^2$ is minimum if and only if $a = \bar{x}$.

(3) If n_1 numbers have mean \bar{x}_1 , n_2 numbers have mean \bar{x}_2 , ..., n_k numbers have mean \bar{x}_k , then the mean of all the numbers called the combined mean is given by

$$\bar{x}_c = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2 + \dots + n_k \bar{x}_k}{n_1 + n_2 + \dots + n_k} = \frac{\sum n_i \bar{x}_i}{\sum n_i}$$

Median

It's the value below which and above which half of the observations fall when ranked in order of size. The position of the median term is given by $\left(\frac{n+1}{2}\right)^{th}$ Value.

NB if n is even we average the middle 2 terms

For grouped data median is estimated using the formula

$$\text{Median} = \text{LCB} + \left(\frac{\left(\frac{n+1}{2}\right) - Cf_a}{f} \right) \times i$$

where LCB, f and i are the lower class boundary, frequency and class interval of the median class. Cf_a is the cumulative frequency of the class above the median class.

Remark: The disadvantage of median is that it is not sensitive against changes in the data.

Mode

It's the value occurring most frequently in a data set. If each observation occurs the same number of times, then there is no mode. When 2 or more observation occurs most frequently in a data then the data is said to be multimodal.

For ungrouped data it's very easy to pick out the mode. However If the data is grouped, mode is estimated using the formula

$$\text{Mode} = \text{LCB} + \left(\frac{f - f_a}{2f - f_a - f_b} \right) \times i$$

where LCB, f and i are the lower class boundary, frequency and class interval of the modal class. f_a and f_b are frequencies of the class above and below the modal class. respectively

Example 1

Find the median and mode of the following data: 19, 13, 18, 14, 12, 25, 11, 10, 17, 23, 19.

Solution

Sorted data: 10, 11, 12, 13, 14, 17, 18, 19, 19, 23, 25.

$n = 11$ thus Median = $\left(\frac{11+1}{2}\right)^{\text{th}}$ Value = 6^{th} Value = 17

Mode=19 since it appears most frequently in this data set as compared to other observations.

Example 2 Find the median and mode of the data: 2, 4, 8, 7, 9, 4, 6, 10, 8, and 5.

Solution

Array: 2, 4, 4, 5, 6, 7, 8, 8, 9, 10

$n = 10$ thus Median = $\left(\frac{10+1}{2}\right)^{\text{th}}$ Value = 5.5^{th} Value = $\frac{6+7}{2} = 6.5$

Mode 4 and 8 ie the data is bimodal.

Example 3

Estimate the mean, median and mode for the following frequency distribution:

Class	5-9	10-14	15-19	20-24	25-29	30-34	35-39
Freq	5	12	32	40	16	9	6

Solution

Boundaries	4.5-9.5	9.5-14.5	14.5-19.5	19.5-24.5	24.5-29.5	29.5-34.5	34.5-39.5
Mid pts (x)	7	12	17	22	27	32	37
Frequency	5	12	32	40	16	9	6
Xf	35	144	544	880	432	288	222
CF	5	17	49	89	105	114	120

$$\text{Mean } \bar{x} = \frac{\sum fx}{n} = \frac{35+144+\dots+222}{120} = \frac{2545}{120} \approx 21.2083$$

$n = 120$ thus Median = 60.5^{th} Value \Rightarrow Median class is 19.5-24.5 thus

$$\text{Median} = \text{LCB} + \left(\frac{\left(\frac{n+1}{2}\right) - Cf_a}{f} \right) \times i = 19.5 + \left(\frac{60.5 - 49}{40} \right) \times 5 = 20.9375$$

The modal class (class with the highest frequency) is 19.5-24.5 therefore

$$\text{Mode} = \text{LCB} + \left(\frac{f - f_a}{2f - f_a - f_b} \right) \times i = 19.5 + 5 \left(\frac{40 - 32}{80 - 32 - 16} \right) = 20.75$$

Exercise

- Find the mean median and mode for the following data: 9, 3, 4, 2, 9, 5, 8, 4, 7, 4
- Find the mean median and mode of 1, 2, 2, 3, 4, 4, 5, 5, 5, 5, 7, 8, 8 and 9
- The number of goals scored in 15 hockey matches is shown in the table.

No of goals	1	2	3	4	5
No of matches	2	1	5	3	4

Calculate the mean number of goals cored

- The table shows the heights of 30 students in a class calculate an estimate of the mean mode and median height.

Height (cm)	140<x<144	144<x<148	148<x<152	152<x<156	156<x<160	160<x<164
No of students	4	5	8	7	5	1

- Estimate the mean, median and mode for the following frequency distribution:

Class	1-4	5-8	9-12	13-16	17-20	21-24
frequency	10	14	20	16	12	8

Class	40-59	60-79	80-99	100-119	120-139	140-159	160-179	180-199
freq	5	12	32	40	16	9	6	

The Empirical Relation between the Mean, Median and Mode

$$MEAN - MODE = 3(MEAN - MEDIAN)$$

The above relation is true for unimodal frequency curves which are asymmetrical.

3.1.2 Other Types of Means

These will include weighted, harmonic and geometric means.

The Weighted Arithmetic Mean

The weighted arithmetic mean of a set of n numbers x_1, x_2, \dots, x_n having corresponding weights w_1, w_2, \dots, w_n is defined as

$$\bar{x}_w = \frac{w_1x_1 + w_2x_2 + \dots + w_nx_n}{w_1 + w_2 + \dots + w_n} = \frac{\sum w_i x_i}{\sum w_i}$$

Example1 Consider the following table with marks obtained by two students James (mark x) and John (mark y). The weights are to be used in determining who joins the engineering course whose requirement is a weighted mean of 58% on the four subjects below;

Subject	Maths	English	History	Physics	Total
Mark x	25	87	83	30	225
Mark y	70	45	35	75	225
Weight	3.6	2.3	1.5	2.6	10

Working the products of the marks and the weights we get

Subject	Maths	English	History	Physics	Total
Wx	90	200.1	124.5	78	492.6
Wy	252	103.5	52.5	195	603

$$\text{Now } \bar{x}_w = \frac{\sum w_i x_i}{\sum w_i} = \frac{492.6}{10} = 49.26 \quad \text{and} \quad \bar{y}_w = \frac{\sum w_i y_i}{\sum w_i} = \frac{603}{10} = 60.3$$

Clearly John qualifies but James does not.

Example 2 If a final examination is weighted 4 times as much as a quiz, a midterm examination 3 times as much as a quiz, and a student has a final examination grade of 80, a midterm examination grade of 95 and quiz grades of 90, 65 and 70, the mean grade is

$$\bar{X} = \frac{1(90) + 1(65) + 1(70) + 3(95) + 4(80)}{1 + 1 + 1 + 3 + 4} = \frac{830}{10} = 83.$$

Question A tycoon has 3 house girls who he pays Ksh 4,000 each per month, 2 watch men who he pays Ksh 5,000 each and some garden men who receives Ksh 7,000 each. If he pays out an average of Ksh 5,700 per month to these people, find the number of garden men.

Question A student's grades in laboratory, lecture, and recitation parts of a computer course were 71, 78, and 89, respectively.

- (a) If the weights accorded these grades are 2, 4, and 5, respectively, what is an average grade?
 (b) What is the average grade if equal weights are used?

The Geometric and Harmonic Means

Let x_1, x_2, \dots, x_n be the sample values, the geometric mean GM is given by

$$GM = \sqrt[n]{x_1 \times x_2 \times \dots \times x_n} = \sqrt[n]{\prod_{i=1}^n x_i}$$

and the harmonic mean is given by

$$HM = \frac{n}{\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_n}} = \frac{n}{\sum \left(\frac{1}{x_i} \right)}$$

The Relation between the Arithmetic, Geometric and Harmonic Means:

$$HM \leq GM \leq \bar{x}.$$

The formulas for geometric and harmonic means of a frequency distribution are respectively given by;

$$GM = \sqrt[n]{x_1^{f_1} \times x_2^{f_2} \times \dots \times x_n^{f_n}} = \sqrt[n]{\prod_{i=1}^n x_i^{f_i}} \Rightarrow \log(GM) = \frac{1}{n} \sum_{i=1}^n f_i \log x_i$$

and

$$HM = \frac{n}{\frac{f_1}{x_1} + \frac{f_2}{x_2} + \dots + \frac{f_n}{x_n}} = \frac{n}{\sum \left(\frac{f_i}{x_i} \right)} = \left[\frac{1}{n} \sum \left(\frac{f_i}{x_i} \right) \right]^{-1}$$

where $n = \sum f_i$ and x_i are the midpoints

Example 1: Find the harmonic and the geometric mean of the numbers 2, 4 and 8

Solution The geometric mean $GM = \sqrt[3]{2 \times 4 \times 8} = \sqrt[3]{64} = 4$ and

the harmonic mean $HM = \frac{3}{\frac{1}{2} + \frac{1}{4} + \frac{1}{8}} = \frac{3}{\frac{7}{8}} = \frac{34}{7} \approx 3.43$

Example 2:

Find the harmonic and geometric mean of the frequency table below

x	13	14	15	16	17
f	2	5	13	7	3

Solution

The harmonic mean $HM = \frac{30}{\frac{2}{13} + \frac{5}{14} + \frac{13}{15} + \frac{7}{16} + \frac{3}{17}} \approx 15$. and

The geometric mean $GM = \sqrt[30]{13^2 \times 14^5 \times 15^{13} \times 16^7 \times 17^3} \approx 15.09837$

Exercise:

- Find the harmonic and the geometric mean of the numbers 10, 12, 15, 5 and 8
- The number of goals scored in 15 hockey matches is shown in the table below . Calculate the harmonic and geometric mean number of goals scored.

No of goals	1	3	5	6	9
No of matches	2	1	5	3	4

- Find the harmonic and geometric mean of the frequency table below

Class	0-29	30-49	50-79	80-99
Frequency	20	30	40	10

3.1.3 Measures of Relative Positioning (Quantiles)

These are values which divide a sorted data set into N equal parts. They are also known as quantiles or N-tiles. The commonly used quantiles are; **Quartiles, Deciles and Percentiles**

These 3 divides a sorted data set into four, ten and hundred divisions, respectively. These measures of position are useful for comparing scores within one set of data. You probably all took some type of college placement exam at some point. If your composite math score was say 28, it might have been reported that this score was in the 94th percentile. What does this mean? This does not mean you received a 94% on the test. It does mean that of all the students who took that exam, 94% of them scored lower than you did (and 6% higher).

Remark For a set of data you can divide the data into three quartiles (Q_1, Q_2, Q_3), nine deciles (D_1, D_2, \dots, D_9) and 99 percentiles (P_1, P_2, \dots, P_{99}). To work with percentiles, deciles and quartiles - you need to learn to do two different tasks. First you should learn how to find the percentile that corresponds to a particular score and then how to find the score in a set of data that corresponds to a given percentile.

Quartiles

They divide a sorted data set into 4 equal parts and we have lower, middle and upper quartiles denoted Q_1 , Q_2 and Q_3 respectively. The lower quartile Q_1 separates the bottom 25% from the top 75%, Q_2 is the median and Q_3 separates the top 25% from the bottom 75% as illustrated below .



The K^{th} quartile is given by: $Q_k = \frac{k}{4}(n+1)^{\text{th}}$ value where $k=1,2,3$

Deciles and Percentiles

Similarly the K^{th} Deciles D_k and the K^{th} Percentiles P_k are respectively given by;

$$D_k = \frac{k}{10}(n+1)^{\text{th}} \text{ Value and } P_k = \frac{k}{100}(n+1)^{\text{th}} \text{ value}$$

NB For ungrouped data we may be forced to use linear interpolation for us to get the required K^{th} quantile. However for grouped data the K^{th} Value is given by

$$K^{\text{th}} \text{ Value} = \text{LCB} + \left(\frac{K - C_{f_a}}{f} \right) \times i$$

where LCB, i and f are the lower class boundary, class interval and frequency of the class containing the K^{th} value. C_{f_a} is the cumulative frequency of the class above this particular class

Example 1: Find the lower and upper quartiles, the 7th decile and the 85th percentile of the following data. 3, 6, 9, 10, 7, 12, 13, 15, 6, 5, 13

Solution

Sorted data: 3, 5, 6, 6, 7, 9, 10, 12, 13, 13, 15 Here $n=11$

$$Q_1 = \frac{1}{4}(11+1)^{\text{th}} = 3^{\text{rd}} \text{ value} = 6 \quad \text{Similarly } Q_3 = \frac{3}{4}(11+1)^{\text{th}} = 9^{\text{th}} \text{ value} = 13$$

$$D_7 = \frac{7}{10}(11+1)^{\text{th}} = 7.7^{\text{th}} \text{ value} = \underbrace{7^{\text{th}} \text{ value} + 0.7(8^{\text{th}} \text{ value} - 7^{\text{th}} \text{ value})}_{\text{linear interpolation}} = 10 + 0.7(12 - 10) = 11.4$$

$$P_{85} = \frac{85}{100}(11+1)^{\text{th}} = 10.2^{\text{th}} \text{ value} = \underbrace{10^{\text{th}} \text{ value} + 0.2(11^{\text{th}} \text{ value} - 10^{\text{th}} \text{ value})}_{\text{linear interpolation}} = 13 + 0.2(15 - 13) = 13.4$$

Example 2:

Estimate the lower quartile, 4th decile and the 72nd percentile for the frequency table below

Class	1-4	5-8	9-12	13-16	17-20	21-24
frequency	10	14	20	16	12	8

Solution

Boundaries	0.5-4.5	4.5-8.5	8.5-12.5	12.5-16.5	16.5-20.5	20.5-24.5
C.F	10	24	44	60	72	80

For this data $n=80$

$$Q_1 = \frac{1}{4}(80+1)^{\text{th}} = 20.25^{\text{th}} \text{ value} = 4.5 + \left(\frac{20.25 - 10}{14} \right) \times 4 \approx 7.428571$$

$$D_4 = \frac{4}{10}(80+1)^{\text{th}} = 32.4^{\text{th}} \text{ value} = 8.5 + \left(\frac{32.4 - 24}{20} \right) \times 4 \approx 10.18$$

$$P_{72} = \frac{72}{100}(80+1)^{\text{th}} = 58.32^{\text{th}} \text{ value} = 12.5 + \left(\frac{58.32 - 44}{16} \right) \times 4 \approx 16.08$$

Exercise

- Find the lower and upper quartiles, the 7th decile and the 85th percentile of the data.
- 9, 3, 4, 2, 9, 5, 8, 4, 7, 4 b) 1, 2, 2, 3, 4, 4, 5, 5, 5, 5, 7, 8, 8 and 9

- 2) The number of goals scored in 15 hockey matches is shown in the table.

No of goals	1	2	3	4	5
No of matches	2	1	5	3	4

Estimate the lower quartile, 4th decile and the 72nd percentile of the number of goals scored

- 4) The table shows the heights of 30 students in a class calculate an estimate of the upper and lower quartile of the height.

Height (cm)	140<x<144	144<x<148	148<x<152	152<x<156	156<x<160	160<x<164
No of students	4	5	8	7	5	1

- 5) The grouped frequency table gives information about the distance each of 150 people travel to work.

Height (cm)	0<d<5	5<d<10	10<d<15	15<d<20	20<d<25	25<d<30
No of students	4	5	8	7	5	1

- a) Work out what percentage of the 150 people travel more than 20 km to work (b) Calculate an estimate for the median distance travelled to work by the people?

Properties of measures of Location

- (i) They are affected by change of origin. Adding or subtracting a constant from each and every observation in a data set causes all the measures of location to shift by the same magnitude. That is $\text{New measure} = \text{old measure} \pm k$
- (ii) They are affected by change of scale. Multiplying each and every observation in a data set by a constant value scales up all the measures of location by the same magnitude.. That is $\text{New measure} = K(\text{old measure})$

Example: Consider the three sets of data A, B and C below

Set A: 65, 53, 42, 52, 53 $\bar{x}_A = 53$ and $\text{Median}_A = 53$

Set B: 15, 3, -8, 2, 3 $\bar{x}_B = 3$ and $\text{Median}_B = 3$

Set C: 45, 9, -24, 6, 9 $\bar{x}_C = 9$ and $\text{Median}_C = 9$

- Notice that set B is obtained by subtracting 50 from each and every observation in set A and clearly $\bar{x}_B = \bar{x}_A - 50$ and $\text{Median}_B = \text{Median}_A - 50$ Therefore $\text{New measure} = \text{old measure} \pm k$. This is referred to as change of origin.
- Effectively set C is obtained by multiplying each and every observation in set B by 3 and clearly $\bar{x}_C = 3\bar{x}_B$ and $\text{Median}_C = 3\text{Median}_B$ Thus $\text{New measure} = K(\text{old measure})$ This is referred to as change of scale.

3.2 Measures of Spread/ Dispersion

Spread is the degree of scatter or variation of the variable about the central value. Examples of these measures includes: the range, Inter-Quartile range, Quartile Deviation also called semi Inter-Quartile range, Mean Absolute Deviation, Variance and standard deviation.

Inter-Quartile range and Semi Inter-Quartile Range

Inter-Quartile range (IQR) is the difference between the upper and lower quartiles. Half of this difference is called Quartile Deviation or the semi Inter-Quartile range (SIQR) Ie

$$\text{IQR} = Q_3 - Q_1 \text{ and } \text{SIQR} = \frac{1}{2}(Q_3 - Q_1)$$

Mean Absolute Deviation (MAD)

It is the average of the absolute deviations from the mean and it's given by

$$\text{MAD} = \frac{\sum |x - \bar{x}|}{n} \text{ for ungrouped data but for grouped data } \text{MAD} = \frac{\sum f |x - \bar{x}|}{n}$$

Example 1:

Find the quartile deviation and the mean absolute deviation for the following data.

3, 6, 9, 10, 7, 12, 13, 15, 6, 5, 13

Solution

Sorted data: 3, 5, 6, 6, 7, 9, 10, 12, 13, 13, 15

Recall $Q_1 = 6$ and $Q_3 = 13$ ie from earlier calculations.

Thus $\text{SIQR} = \frac{1}{2}(Q_3 - Q_1) = \frac{1}{2}(13 - 6) = 3.5$

$$\bar{x} = \frac{3+5+6+6+7+9+10+12+13+13+15}{11} = 9$$

$$\text{MAD} = \frac{\sum (x - \bar{x})}{n} = \frac{|3-9| + |5-9| + |6-9| + \dots + |15-9|}{11} = \frac{6+4+3+\dots+6}{11} = \frac{36}{11} \approx 3.2727$$

Variance and Standard Deviation

Ignoring the negative sign in order to compute MAD is not the only option we have to deal with deviations. We can square the deviations and then average. The average of the squared deviations from the mean is called the variance denoted s^2 and its given by

$$s^2 = \frac{1}{n} \sum (x - \bar{x})^2 \text{ A little algebraic simplification of this formular gives } s^2 = \frac{1}{n} \sum x^2 - \bar{x}^2$$

For grouped data $s^2 = \frac{1}{n} \sum f(x - \bar{x})^2 = \frac{1}{n} \sum fx^2 - \bar{x}^2$ where n is the sum of frequencies.

To reverse the squaring on the units we find the square root of the variance. Standard Deviation denoted s is the square root of variance.

Example 1: Find the variance and standard deviation for the data.

3, 6, 9, 10, 7, 12, 13, 15, 6, 5, 13

Solution

$$\bar{x} = \frac{3+5+6+6+7+9+10+12+13+13+15}{11} = 9$$

$$s^2 = \frac{\sum (x - \bar{x})^2}{n} = \frac{(3-9)^2 + (5-9)^2 + (6-9)^2 + \dots + (15-9)^2}{11} = \frac{36+16+9+\dots+36}{11} = \frac{143}{11} = 13$$

$$\text{Standard deviation } s = \sqrt{\text{variance}} = \sqrt{13} \approx 3.60555.$$

Example 2 Find the standard deviation of the data: 2, 4, 8, 7, 9, 4, 6, 10, 8, and 5.

Solution

$$\text{Mean } \bar{x} = \frac{\sum x}{n} = \frac{2+4+8+\dots+5}{10} = \frac{63}{10} = 6.3 \text{ and } \sum x^2 = 2^2 + 4^2 + 8^2 + \dots + 5^2 = 455$$

$$\text{Standard deviation } s = \sqrt{\frac{1}{n} \sum x^2 - \bar{x}^2} = \sqrt{45.5 - 6.3^2} \approx 2.4104.$$

Example 3 Estimate the mean, and standard deviation for the frequency table below:

Class	5-9	10-14	15-19	20-24	25-29	30-34	35-39
freq	5	12	32	40	16	9	6

Solution

Mid pts (x)	7	12	17	22	27	32	37	sum
Freq (f)	5	12	32	40	16	9	6	120
xf	35	144	544	880	432	288	222	2545
fx^2	245	1728	9248	19360	11664	9216	8214	59675

$$\text{Mean } \bar{x} = \frac{\sum fx}{n} = \frac{2545}{120} \approx 21.2083 \text{ and } \sum fx^2 = 59675$$

$$\text{Standard deviation } s = \sqrt{\frac{1}{n} \sum fx^2 - \bar{x}^2} = \sqrt{\frac{59675}{120} - 21.2083^2} \approx 6.8919.$$

Exercise

- Find the quartile deviation, the mean absolute deviation and the standard deviation of the data: a) 9, 3, 4, 2, 9, 5, 8, 4, 7, 4 b) 1, 2, 2, 3, 4, 4, 5, 5, 5, 5, 7, 8, 8 and 9
- The number of goals scored in 20 hockey matches is shown in the table.

No of goals	1	2	3	4	5
No of matches	2	5	6	3	4

Estimate the quartile deviation, the mean absolute deviation and the standard deviation of the number of goals scored

- consider the frequency table below and estimate quartile deviation, the mean absolute deviation and the standard deviation

Class	8-12	13-17	18-22	23-27	28-32	33-37
Freq	3	10	12	9	5	1

- The table shows the heights of 30 students in a class calculate an estimate of the quartile deviation, the mean absolute deviation and the standard deviation of the height.

Height (cm)	140<x<144	144<x<148	148<x<152	152<x<156	156<x<160	160<x<164
No of students	4	5	8	7	5	1

- The grouped frequency table gives information about the distance each of 150 people travel to work.

Height (cm)	0<d<5	5<d<10	10<d<15	15<d<20	20<d<25	25<d<30
No of students	4	5	8	7	5	1

Calculate an estimate for the quartile deviation and the standard deviation of the distance travelled to work by the people.

Properties of measures of Spread

- They are not affected by change of origin. Adding or subtracting a constant from each and every observation in a data set does not affect any measures of spread. That is New measure = old measure
- They are affected by change of scale. Multiplying each and every observation in a data set by a constant value scales up all the measures of spread by the

same value except in the case of variance which is scaled up by a square of the same constant.

ie New measure = $K(\text{old measure})$ but New variance = $k^2 \times \text{old variance}$

Example: Consider the three sets of data A, B and C below

Set A: 65, 53, 42, 52, 53 Range=23, $MAD_A = 4.8$ and $Variance_A = 66.5$

Set B: 15, 3, -8, 2, 3 Range=23, $MAD_B = 4.8$ and $Variance_B = 66.5$

Set C: 45, 9, -24, 6, 9 Range=69, $MAD_C = 14.4$ and $Variance_C = 598.5$

- Notice that set B is obtained by subtracting 50 from each and every observation in set A and clearly $MAD_B = MAD_A$ and $Variance_B = Variance_A$ Therefore there is no effect on the change of origin ie New measure = old measure..
- Effectively set C is obtained by multiplying each and every observation in set B by 3 and clearly $MAD_C = 3 \times MAD_B$ and $Variance_C = 3^2 \times Variance_B$ Thus
New measure = $K(\text{old measure})$ and New $Variance_C = k^2 \times \text{old } Variance_B$

Mean and Standard Deviation Using a Calculator

- When on, press mode key to get;
COMP SD REG
1 2 3
- Press 2 to select SD for statistical data.
- Enter data one by one pressing m+ after every value entered. The screen will be showing the number of observations that are fully entered.
- Pressing shift then 1 gives

$$\begin{array}{ccc} \sum x^2 & \sum x & n \\ 1 & 2 & 3 \end{array}$$
- Pressing shift then 2 gives

$$\begin{array}{ccc} \bar{x} & x\sigma_n & x\sigma_{n-1} \\ 1 & 2 & 3 \end{array}$$

Typing 1 then = gives the value of $\sum x^2$
Similarly typing 2 then = gives the value of $\sum x$

Which are the mean uncorrected standard deviation and the corrected standard deviation

Example using your calculator, obtain the mean and standard deviation of the following data: 31, 52, 29, 60, 58

Solution

Entering data 31M+ 52 M+ 29 M+ 60 M+ 58 M+

$\bar{x} = 46$ and $s = 14.91643$

Question Redo the above example using the data: 235, 693, 484, 118, 470

3.3 Assumed Mean and the Coding Formular

If the observations are too large such that the natural computation of totals is tedious, we can take one of the observations as the working/assumed mean. Let A be any guessed or assumed arithmetic mean and let $d_i = x_i - A$ be the deviations of x_i from A, then

$$\text{mean } \bar{x} = A + \frac{1}{n} \sum fd = A + \bar{d}$$

and

$$\text{Variance } S^2 = \frac{1}{n} \sum f d^2 - \left(\frac{1}{n} \sum f d \right)^2 = \frac{1}{n} \sum f d^2 - \bar{d}^2$$

Respectably where A = Assumed mean which is generally taken as mid point of the middle class or the class where frequency is large

Remark:, in most cases deviations (d) of x_i from A is a multiple of the class interval ie

$$d_i = t_i \times i \Rightarrow t_i = \frac{d_i}{i} = \frac{x_i - A}{i}.$$

In these cases we can use t rather than d in computation. The above formulae reduces to

$$\bar{x} = A + \frac{i}{n} \sum f t = A + i \bar{t} \quad \text{and}$$

$$S^2 = i^2 \left[\frac{1}{n} \sum f t^2 - \left(\frac{1}{n} \sum f t \right)^2 \right] = i^2 \left[\frac{1}{n} \sum f t^2 - \bar{t}^2 \right]$$

respectably the latter formulae are referred to as coding formulae

Example

Using coding formulae, find the mean and standard deviation of the following data

Class	340-349	350-359	360-369	370-379	380-389
Freq	2	3	7	5	3

Solution

Class	Mid pts	Freq	$t = \frac{x - 364.5}{i}$	ft	ft ²
340-349	344.5	2	-2	-4	8
350-359	354.5	3	-1	-3	3
360-369	364.5	7	0	0	0
370-379	374.5	5	1	5	5
380-389	384.5	3	2	6	12
Total		20		4	28

$$\bar{x} = A + \frac{i}{n} \sum f t = 364.5 + \frac{10}{20} (4) = 366.5$$

$$S = i \times \sqrt{\frac{1}{n} \sum f t^2 - \left(\frac{1}{n} \sum f t \right)^2} = 10 \times \sqrt{\frac{28}{20} - \left(\frac{4}{20} \right)^2} \approx 11.6619$$

Exercise

- 1) Consider the following frequency distribution.

classes	10-14	15-19	20-24	25-29	30-34
frequency	7	11	14	13	5

Estimate the mean and standard deviation using coding formula

- 2) Using coding formular, find the mean and standard deviation of the frequency table below

Class	10-20-	20-30	30-40	40-50	50-60	60-70	70-80	80-90	90-100	100-110
Freq	4	5	7	13	16	11	9	6	4	3

- 3) The table shows the speed distribution of vehicles on Thika Super high way on a typical day.

Speed (km/hr)	60-69	70-79	80-89	90-99	100-109	110-119	120-129	130-139	140-149
No of vehicles	138	163	325	541	427	214	110	52	30

Using coding formulae, find the mean speed and the standard deviation of the speeds.

- 4) The following table shows a frequency distribution of the weekly wages of 65 employees at the P&R Company.

Wages	\$250.00-259.99	\$260.00-269.99	\$270.00-279.99	\$280.00-289.99	\$290.00-299.99	\$300.00-309.99	\$310.00-319.99
No. of employees	8	10	16	14	10	5	2

Find the mean wage and the standard deviation of the wages using coding formular

3.4 Measures of Relative Dispersion:

These measures are used in comparing spreads of two or more sets of observations. These measures are independent of the units of measurement. These are a sort of ratio and are called coefficients.

Suppose that the two distributions to be compared are expressed in the same units and their means are equal or nearly equal. Then their variability can be compared directly by using their standard deviations. However, if their means are widely different or if they are expressed in different units of measurement, we can not use the standard deviations as such for comparing their variability. We have to use the relative measures of dispersion in such situations. Examples of these Measures of relative dispersion includes; Coefficient of quartile deviation, Coefficient of mean deviation and the Coefficient of variation

3.4.1 Coefficient of Quartile Deviation and Coefficient of Mean Deviation

The Coefficient of Quartile Deviation of x CQD(x) is given by $CQD(x) = \frac{Q_3 - Q_1}{Q_3 + Q_1} \times 100\%$

The Coefficient of Mean Deviation CMD(x) is given by $CMD(x) = \frac{MAD}{Mean} \times 100\%$

3.4.2 Coefficient of Variation:

Coefficient of variation is the percentage ratio of standard deviation and the arithmetic mean. It is usually expressed in percentage. The coefficient of variation of x denoted C.V(x) is given by the formula

$$C.V(x) = \frac{s}{\bar{x}} \times 100\%$$

where \bar{x} is the mean and S is the standard deviation of x.

The coefficient has no units ie it's independent of the units of measurements. It is useful in comparing spreads of two or more populations. The smaller the coefficient of variation, the higher the peak and the lower the spread.

Note: Standard deviation is absolute measure of dispersion while. Coefficient of variation is relative measure of dispersion.

Example 1: Consider the distribution of the yields (per plot) of two ground nut varieties. For the first variety, the mean and standard deviation are 82 kg and 16 kg respectively. For the second variety, the mean and standard deviation are 55 kg and 8 kg respectively. Then we have, for the first variety

$$C.V(x) = \frac{16}{82} \times 100 \approx 19.5\%$$

For the second variety

$$C.V(x) = \frac{8}{55} \times 100 \approx 14.5\%$$

It is apparent that the variability in second variety is less as compared to that in the first variety. But in terms of standard deviation the interpretation could be reverse.

Example 2: Below are the scores of two cricketers in 10 innings. Find who is more „consistent scorer“ by Indirect method.

A	204	68	150	30	70	95	60	76	24	19
B	99	190	130	94	80	89	69	85	65	40

Solution:

From a calculator, $\bar{x}_A = 79.6$, $S_A = 58.2$ $\bar{x}_B = 94.1$ and $S_B = 41.1$

Coefficient of variation for player A is $C.V(x) = \frac{58.2}{79.6} \times 100 \approx 73.153\%$

Coefficient of variation for player B is $C.V(x) = \frac{41.1}{94.1} \times 100 \approx 43.7028\%$

Coefficient of variation of A is greater than coefficient of variation of B and hence we conclude that player B is more consistent

Exercise

- Find the coefficient of quartile deviation, the coefficient of mean deviation and the Coefficient of variation n of x for the following data:
 - 9, 3, 4, 2, 9, 5, 8, 4, 7, 4
 - 1, 2, 2, 3, 4, 4, 5, 5, 6, 6, 7, 8, 8 and 9
 - 3, 6, 9, 10, 7, 12, 13, 15, 6, 5, 13
 - data on marks given by the table below

Marks Obtained	0-10	10-20	20-30	30-40	40-50	50-60	60-70
No. of Students	6	12	22	24	16	12	8

- If the weights of 7 ear-heads of sorghum are 89, 94, 102, 107, 108, 115 and 126 g. Find the arithmetic mean and standard deviation using a calculator hence determine the coefficient of variation of the ear-heads of sorghum
- The following are the 381 soybean plant heights in Cms collected from a particular plot. Using coding formula, Find the mean and Standard deviation of the plants hence determine the coefficient of variation of the 1 soybean plant heights:

Plant heights (Cms)	6.8-7.2	7.3-7.7	7.8-8.2	8.3-8.7	8.8-9.2	9.3-9.7	9.8-10.2	10.3-10.7	10.8-11.2	11.3-11.7	11.8-12.2	12.3-12.7
No. of Plants	9	10	11	32	42	58	65	55	37	31	24	7

3.5 Measures of Skewness and Kurtosis

3.5.1 Skewness

Before discussing the concept of skewness, an understanding of the concept of **symmetry** is essential. A plot of frequency against class mark joined with a smooth curve can help us to visually assess the symmetry of a distribution. Usually symmetry is about the central value. Symmetry is said to exist in a distribution if the smoothed frequency polygon of the distribution can be divided into two identical halves wherein each half is a mirror image of the other. **Skewness** on the other hand means lack of symmetry and it can be positive or negative.

Basically, if the distribution has a tail on the right, (See figure below), then the distribution is positively skewed. Eg Most students having very low marks in an examination. However if the distribution has a tail on the left, then the distribution is negatively skewed. (see figure below). Eg Most students having very high marks in an

examination

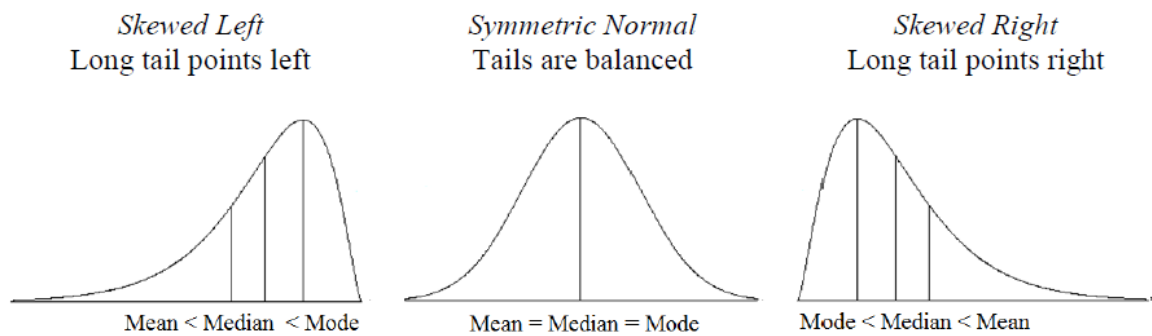


Figure 1. Sketches showing general position of mean, median, and mode in a population.

Measures of Skewness

Generally for any set of values $x_1, x_2, x_3, \dots, x_n$, the **moment coefficient of skewness** α_3 is

given by $\alpha_3 = \frac{\sum f(x - \bar{x})^3}{nS^3}$ where S is the standard deviation of X . It's worth noting that if $\alpha_3 < 0$, the distribution is negatively skewed, if $\alpha_3 > 0$, the distribution is positively skewed and if $\alpha_3 = 0$ the distribution is normal

Other measures of Skewness includes the Karl Pearson coefficient of Skewness SK_p ,

Bowley's coefficient of Skewness SK_B and Kelley's coefficient of Skewness SK_k .

The **Karl Pearson's coefficient of Skewness** is based upon the **divergence of mean from mode** in a skewed distribution. Recall the empirical relation between mean, median and mode which states that, for a moderately symmetrical distribution, we have

$$\text{Mean} - \text{Mode} = 3(\text{Mean} - \text{Median})$$

Hence Karl Pearson's coefficient of skewness is defined by;

$$SK_p = \frac{\text{Mean} - \text{Mode}}{\text{Standard Deviation}} = \frac{3(\text{Mean} - \text{Median})}{\text{Standard Deviation}}$$

The **Bowley's coefficient of Skewness** is based on quartiles. For a symmetrical distribution, it is seen that Q_1 and Q_3 **are** equidistant from median.

$$SK_B = \frac{Q_3 - 2Q_2 + Q_1}{Q_3 - Q_1} \text{ where } Q_k \text{ is the } K^{\text{th}} \text{ quartile.}$$

The **Kelly's coefficient of Skewness** is based on P_{90} and P_{10} so that only 10% of the observations on each extreme are ignored.. This is an improvement over the Bowley's coefficient which leaves 25% of the observations on each extreme of the distribution.

$$SK_k = \frac{P_{90} - 2P_{50} + P_{10}}{P_{90} - P_{10}} \text{ where } P_k \text{ is the } K^{\text{th}} \text{ percentile.}$$

Interpreting Skewness

If the coefficient of skewness is positive, the data are positively skewed or skewed right, meaning that the right tail of the distribution is longer than the left. If the coefficient of skewness is negative, the data are negatively skewed or skewed left, meaning that the left tail is longer. If the coefficient of skewness = 0, the data are perfectly symmetrical. But a skewness of exactly zero is quite unlikely for real-world data, so *how can you interpret the*

skewness number? Bulmer, M. G., *Principles of Statistics* (Dover, 1979) — a classic — suggests this rule of thumb: If the coefficient of skewness is:-

- less than -1 or greater than $+1$, the distribution is *highly skewed*.
- between -1 and $-\frac{1}{2}$ or between $+\frac{1}{2}$ and $+1$, the distribution is *moderately skewed*.
- between $-\frac{1}{2}$ and $+\frac{1}{2}$..., the distribution is *approximately symmetric*.

Example: The following figures relate to the size of capital of 285 companies :

Capital (in Ks lacs.)	1-5	6-10	11-15	16-20	21-25	26-30	31-35
No. of companies	20	27	29	38	48	53	70

Compute the Bowley's coefficients of skewness and interpret the results.

Solution

Boundaries	0.5-5.5	5.5-10.5	10.5-15.5	15.5-20.5	20.5-25.5	25.5-30.5	30.5-35.5
CF	20	47	76	114	162	215	285

$$Q_1 = \frac{1}{4}(286)^{\text{th}} \text{ value} = 71.5^{\text{th}} \text{ value} = 10.5 + \left(\frac{71.5 - 47}{29} \right) \times 5 \approx 14.7241$$

$$Q_2 = \frac{1}{2}(286)^{\text{th}} \text{ value} = 143^{\text{rd}} \text{ value} = 20.5 + \left(\frac{143 - 114}{48} \right) \times 5 \approx 23.5208$$

$$Q_3 = \frac{3}{4}(286)^{\text{th}} \text{ value} = 214.5^{\text{th}} \text{ value} = 25.5 + \left(\frac{214.5 - 162}{53} \right) \times 5 \approx 30.4528$$

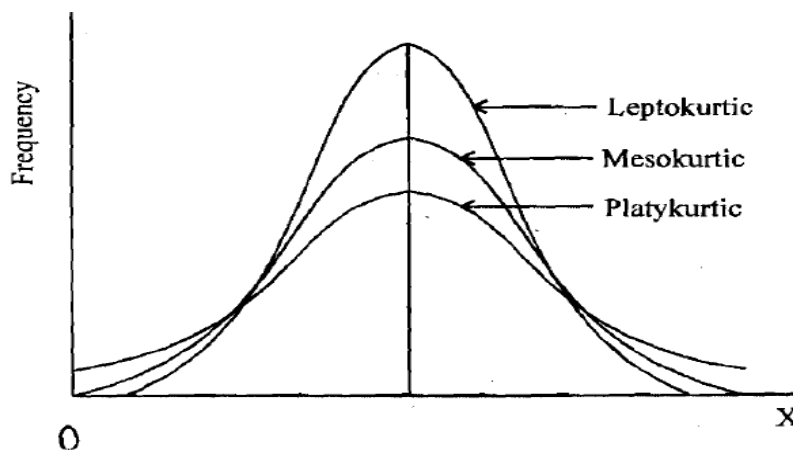
$$SK_p = \frac{Q_3 - 2Q_2 + Q_1}{Q_3 - Q_1} = \frac{30.4528 - 2 \times 23.5208 + 14.7241}{30.4528 - 14.7241} \approx -0.11855.$$

This value lies between $-\frac{1}{2}$ and $+\frac{1}{2}$, therefore the distribution is approximately symmetric.

Question: Compute the Karl Pearson's and the Kelly's coefficient of skewness for the above data and interpret the results.

3.5.2 Kurtosis

It measures the peakedness of a distribution. If the values of x are very close to the mean, the peak is very high and the distribution is said to be **Leptokurtic**. On the other hand if the values of x are very far away from the mean, the peak is very low and the distribution is said to be **Platykurtic**. Finally if x values are at a moderate distance from the mean then the peak is moderate and the distribution is said to be **mesokurtic**.



Measures of Kurtosis

Generally for a set of values $x_1, x_2, x_3, \dots, x_n$, the moment coefficient of kurtosis α_4 is given

by $\alpha_4 = \frac{\sum f(x - \bar{x})^4}{nS^4}$ where \bar{x} and S are the arithmetic mean and standard deviation of X.

Example: Calculate the coefficient of Skewness α_3 and the coefficient of kurtosis α_4 for the data 5, 6, 7, 6, 9, 4, 5

Solution

$$\bar{x} = \frac{1}{n} \sum x = \frac{42}{7} = 6 \quad \text{and} \quad \text{Standard deviation } s = \sqrt{\frac{1}{n} \sum (x - \bar{x})^2} = \frac{4}{\sqrt{7}}$$

x	5	6	7	6	9	4	5	Sum
$(x - \bar{x})^2$	1	0	1	0	9	4	1	16
$(x - \bar{x})^3$	-1	0	1	0	27	-8	-1	18
$(x - \bar{x})^4$	1	0	1	0	81	16	1	100

$$\text{Coefficient of Skewness } \alpha_3 = \frac{\sum (x - \bar{x})^3}{nS^3} = \frac{18}{7} \times \left(\frac{\sqrt{7}}{4}\right)^3 \approx 0.744118$$

$$\text{Coefficient of kurtosis } \alpha_4 = \frac{\sum f(x - \bar{x})^4}{nS^4} = \frac{100}{7} \times \left(\frac{\sqrt{7}}{4}\right)^4 \approx 2.73438$$

Exercise

- Find the moment coefficient of Skewness and kurtosis for the data below. a) 9, 3, 4, 2, 9, 5, 8, 4, 7, 4 b) 1, 2, 2, 3, 4, 4, 5, 5, 6, 6, 7, 8, 8 and 9 c) 3, 6, 9, 10, 7, 12, 13, 15, 6, 5, 13
d) data on marks given by the table below

Marks Obtained	0-10	10-20	20-30	30-40	40-50	50-60	60-70
No. of Students	6	12	22	24	16	12	8

- Data given by the table below

Marks Obtained	0-10	10-20	20-30	30-40
No. of Students	1	3	4	2

- Compute the Bowley's coefficient of skewness, the Kelly's coefficient of skewness and the Percentile coefficient of kurtosis for the following data and interpret the results.

- 9, 3, 4, 2, 9, 5, 8, 4, 7, 4 b) 1, 2, 2, 3, 4, 4, 5, 5, 6, 6, 7, 8, 8 and 9
c) 3, 6, 9, 10, 7, 12, 13, 15, 6, 5, 13 d) data on heights given by the table below

Height (in inches.)	58	59	60	61	62	63	64	65
No. of persons	10	18	30	42	35	28	16	8

- data on daily expenditure of families given by the table below

Daily Expenditure (Rs)	0-20	20-40	40-60	60-80	80-100
No. of persons	13	25	27	19	16

- Data on marks given by the table below

Marks Obtained	0-20	20-40	40-60	60-80	80-100
No. of Students	8	28	35	17	12

- The following measures were computed for a frequency distribution :
Mean = 50, coefficient of Variation = 35% and Karl Pearson's Coefficient of Skewness $SK_p = -0.25$. Compute Standard Deviation, Mode and Median of the distribution.

4. Bivariate Data

4.1 Introduction

So far we have confined our discussion to the distributions involving only one variable. Sometimes, in practical applications, we might come across certain set of data, where each item of the set may comprise of the values of two or more variables.

A Bivariate Data is a set of paired measurements which are of the form

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

Examples

- i. Marks obtained in two subjects by 60 students in a class.
- ii. The series of sales revenue and advertising expenditure of the various branches of a company in a particular year.
- iii. The series of ages of husbands and wives in a sample of selected married couples.

In a bivariate data, each pair represents the values of the two variables. Our interest is to find a relationship (if it exists) between the two variables under study.

4.2 Scatter Diagrams and Correlation

A scatter diagram is a tool for analyzing relationships between two variables. One variable is plotted on the horizontal axis and the other is plotted on the vertical axis. The pattern of their intersecting points can graphically show relationship patterns. Most often a scatter diagram is used to prove or disprove cause-and-effect relationships. While the diagram shows relationships, it does not by itself prove that one variable *causes* the other. In brief, the easiest way to visualize Bivariate Data is through a Scatter Plot.

“Two variables are said to be correlated if the change in one of the variables results in a change in the other variable”.

4.2.1: Positive and Negative Correlation

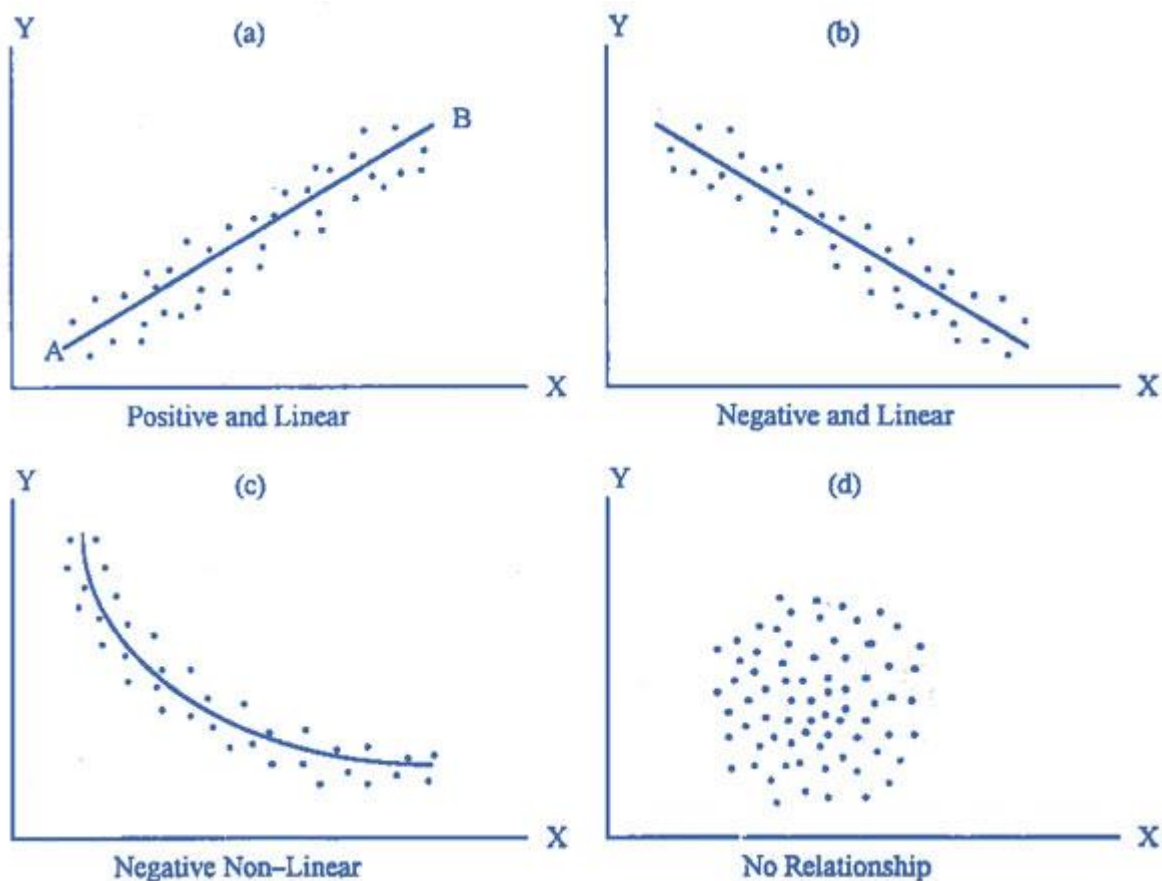
If the values of the two variables deviate in the same direction i.e. if an increase (or decrease) in the values of one variable results, on an average, in a corresponding increase (or decrease) in the values of the other variable the correlation is said to be positive.

Some examples of series of positive correlation are:

- i. Heights and weights;
- ii. Household income and expenditure;
- iii. Price and supply of commodities;
- iv. Amount of rainfall and yield of crops.

Correlation between two variables is said to be negative or inverse if the variables deviate in opposite direction. That is, if the increase in the variables deviate in opposite direction. That is, if increase (or decrease) in the values of one variable results on an average, in corresponding decrease (or increase) in the values of other variable.

Eg Price and demand of goods.



4.2.2 Interpreting a Scatter Plot

Scatter diagrams will generally show one of six possible correlations between the variables:

- i. *Strong Positive Correlation* The value of Y clearly increases as the value of X increases.
- ii. *Strong Negative Correlation* The value of Y clearly decreases as the value of X increases.

- iii. *Weak Positive Correlation* The value of Y increases slightly as the value of X increases.
- iv. *Weak Negative Correlation* The value of Y decreases slightly as the value of X increases.
- v. *Complex Correlation* The value of Y seems to be related to the value of X, but the relationship is not easily determined.
- vi. *No Correlation* There is no demonstrated connection between the two variables

4.3 Correlation Coefficient

Correlation coefficient measures the degree of linear association between 2 paired variables. It takes values from +1 to -1.

- i. If $r = +1$, we have **perfect positive** relationship
- ii. If $r = -1$, we have **perfect negative** relationship
- iii. If $r = 0$ there is **no** relationship i.e. the variables are **uncorrelated**.

4.3.1 Pearson's Product Moment Correlation Coefficient

Pearson's product moment correlation coefficient, usually denoted by r , is one example of a correlation coefficient. It is a measure of the linear association between two variables that have been measured on interval or ratio scales, such as the relationship between height in inches and weight in pounds. However, it can be misleadingly small when there is a relationship between the variables but it is a non-linear one.

The correlation coefficient r is given by
$$r = \frac{n \sum xy - \sum x \sum y}{\sqrt{[n \sum x^2 - (\sum x)^2][n \sum y^2 - (\sum y)^2]}}$$

Example:

: A study was conducted to find whether there is any relationship between the weight and blood pressure of an individual. The following set of data was arrived at from a clinical study. Let us determine the coefficient of correlation for this set of data. The first column represents the serial number and the second and third columns represent the weight and blood pressure of each patient.

Weight	78	86	72	822	80	86	84	89	68	71
Blood Pressure	140	160	134	144	180	176	174	178	128	132

Solution:

x	y	x ²	y ²	xy
78	140	6084	19600	10920
86	160	7396	25600	13760
72	134	5184	17956	9648
82	144	6724	20736	11808
80	180	6400	32400	14400
86	176	7396	30976	15136
84	174	7056	30276	14616
89	178	7921	31684	15842
68	128	4624	16384	8704
71	132	5041	17424	9372
796	1546	63,776	243036	1242069

Thus

$$r = \frac{10(124206) - (796)(1546)}{\sqrt{[(10)63776 - (796)^2](10)[(243036) - (1546)^2]}} = \frac{11444}{\sqrt{(1144)(40244)}} = 0.5966$$

It can be shown that $r = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2} \sqrt{\sum (y - \bar{y})^2}}$

Example:

Obtain the correlation coefficient of the following data

Mean Temp. (x)	14.2	14.3	14.6	14.9	15.2	15.6	15.9
Pirates (y)	35000	45000	20000	15000	5000	400	17

Solution

Mean Temp. (x)	Pirates (y)	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})^2$	$(y_i - \bar{y})^2$	$(x_i - \bar{x})(y_i - \bar{y})$
14.2	35000	-0.76	17797.57	0.57	316753548	-13475
14.3	45000	-0.66	27797.57	0.43	772704977	-18266
14.6	20000	-0.36	2797.57	0.13	7826405	-999
14.9	15000	-0.06	-2202.43	0	4850691	125
15.2	5000	0.24	-12202.43	0.06	148899263	-2963
15.6	400	0.64	-16802.43	0.41	282321605	-10801
15.9	17	0.94	-17185.43	0.89	295338955	-16203
Tot.=104.7	120417	0	0	2.5	1828695447	-62583
$\bar{x} = 14.96$	$\bar{y} = 17202.43$			S_{xx}	S_{yy}	S_{xy}

We then have that $r = \frac{-62583}{\sqrt{2.5(1828695447)}} \approx -0.93$

4.3.2 Spearman rank correlation coefficient

Data which are arranged in ascending order are said to be in **ranks** or **ranked data**. The coefficient of correlation for such type of data is given by **Spearman rank difference correlation coefficient** and is denoted by R.

R is given by the formula $R = 1 - \frac{6\sum d^2}{n(n^2 - 1)}$

Example

The data given below are obtained from student records. (Grade Point Average (x) and Graduate Record exam score (y)) Calculate the rank correlation coefficient 'R' for the data.

Subject	1	2	3	4	5	6	7	8	9	10
X	8.3	8.6	9.2	9.8	8.0	7.8	9.4	9.0	7.2	8.6
y	2300	2250	2380	2400	2000	2100	2360	2350	2000	2260

Solution

Note that in the x row, we have two students having a grade point average of 8.6 also in the y row; there is a tie for 2000.

Now we arrange the data in descending order and then rank 1,2,3,.. . .10 accordingly. In case of a tie, the rank of each tied value is the mean of all positions they occupy. In x, for

instance, 8.6 occupy ranks 5 and 6. So each has a rank $\frac{5+6}{2} = 5.5$

Similarly in 'y' 2000 occupies ranks 9 and 10, so each has rank 9.5

Now we come back to our formula $R = 1 - \frac{6\sum d^2}{n(n^2 - 1)}$

We compute d, square it and substitute its value in the formula

Subject	x	y	Rank of x	Rank of y	d	d ²
1.	8.3	2300	7	5	2	4
2.	8.6	2250	5.5	7	-1.5	2.25
3.	9.2	2380	3	2	1	1
4.	9.8	2400	1	1	0	0
5.	8.0	2000	8	9.5	-1.5	2.25
6.	7.8	2100	9	8	1	1
7.	9.4	2360	2	3	-1	1
8.	9.0	2350	4	4	0	0
9.	7.2	2000	10	9.5	0.5	0.25
10.	8.6	2260	5.5	6	-0.5	0.25

So here, n = 10 and $\sum d^2 = 12$. So

$$R = 1 - \frac{6(12)}{10(100 - 1)} = 1 - 0.0727 = 0.9273$$

Note: If we are provided with only ranks without giving the values of x and y we can still find Spearman rank difference correlation R by taking the difference of the ranks and proceeding in the above shown manner.

4.4 Regression

If two variables are significantly correlated, and if there is some theoretical basis for doing so, it is possible to predict values of one variable from the other.

Regression analysis, in general sense, means the estimation or prediction of the unknown value of one variable from the known value of the other variable. It is one of the

most important statistical tools which is extensively used in almost all sciences – Natural, Social and Physical.

Regression analysis was explained by M. M. Blair as follows:

“Regression analysis is a mathematical measure of the average relationship between two or more variables in terms of the original units of the data.”

3.4.1 Regression Equation

Regression analysis can be thought of as being sort of like the flip side of correlation. It has to do with finding the equation for the kind of straight lines you were just looking at. Suppose we have a sample of size n and it has two sets of measures, denoted by x and y . We can predict the values of y given the values of x by using the equation, $y^* = a + bx$

Where the coefficients ‘ a ’ and ‘ b ’ are real numbers given by

$$b = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2} \quad \text{and} \quad a = \frac{\sum y - b \sum x}{n}$$

The symbol y^* refers to the predicted value of y from a given value of x from the regression equation.

Example:

Scores made by students in a statistics class in the mid-term and final examination are given here. Develop a regression equation which may be used to predict final examination scores from the mid – term score.

Student	1	2	3	4	5	6	7	8	9	10
Mid term	98	66	100	96	88	45	76	60	74	82
Final	90	74	98	88	80	62	78	74	86	80

Solution:

We want to predict the final exam scores from the mid term scores. So let us designate ‘ y ’ for the final exam scores and ‘ x ’ for the mid term exam scores. We open the following table for the calculations.

Stud	x	y	X^2	xy
1	98	90	9604	8820
2	66	74	4356	4884
3	100	98	10,000	9800
4	96	88	9216	8448
5	88	80	7744	7040
6	45	62	2025	2790
7	76	78	5776	5928
8	60	74	3600	4440
9	74	86	5476	6364
10	82	80	6724	6560
Total	785	810	64,521	65,071

$$b = \frac{10(65,071) - 785(810)}{10(64,521) - (785)^2} = \frac{14,860}{28,985} = 0.5127 \quad \text{and} \quad a = \frac{810 - 785(0.5127)}{10} = 40.7531$$

Thus, the regression equation is given by $y^* = 40.7531 + (0.5127)x$

We can use this to find the projected or estimated final scores of the students.

Eg for the midterm score of 50 the projected final score is

$$y^* = 40.7531 + (0.5127)50 = 66.3881, \text{ which is a quite a good estimation.}$$

To give another example, consider the midterm score of 70. Then the projected final score is

$$y^* = 40.7531 + (0.5127)70 = 76.6421, \text{ which is again a very good estimation.}$$

Practice Problems:

1. Consider the following data and draw a scatter plot

X	1.0	1.9	2.0	2.9	3.0	3.1	4.0	4.1	5
Y	10	99	100	999	1,000	1,001	10,000	10,001	100,000

2. Let variable X is the number of hamburgers consumed at a cook-out, and variable Y is the number of beers consumed. Develop a regression equation to predict how many beers a person will consume given that we know how many hamburgers that person will consume.

Subject	1	2	3	4	5
Hamburgers	5	4	3	2	1
Beers	8	10	4	6	2

3. A horse owner is investigating the relationship between weight carried and the finish position of several horses in his stable. Calculate r and R for the data given

Weight carried	110	113	120	115	110	115	117	123	106	108	110	110
Position Finished	2	6	3	4	6	5	4	2	1	4	1	3

4. The top and bottom number which may appear on a die are as follows Calculate r and R for these values. Are the results surprising?

Top	1	2	3	4	5	6
Bottom	5	6	4	3	1	2

5. The ranks of two sets of variables (Heights and Weights) are given below. Calculate the Spearman rank difference correlation coefficient R.

	1	2	3	4	5	6	7	8	9	10
Heights	2	6	8	4	7	4	9.5	4	1	9.5
Weights	9	1	9	4	5	9	2	7	6	3

6. Researchers interested in determining if there is a relationship between death anxiety and religiosity conducted the following study. Subjects completed a death anxiety scale (high score = high anxiety) and also completed a checklist designed to measure an individuals degree of religiosity (belief in a particular religion, regular attendance at religious services, number of times per week they regularly pray, etc.) (high score = greater religiosity). A data sample is provided below:

X	38	42	29	31	28	15	24	17	19	11	8	19	3	14	6
y	4	3	11	5	9	6	14	9	10	15	19	17	10	14	18

- What is your computed answer?
- What does this statistic mean concerning the relationship between death anxiety and religiosity?
- What percent of the variability is accounted for by the relation of these two variables?

7. The data given below are obtained from student records.(Grade Point Average (x) and Graduate Record exam score (y)) Calculate the regression equation and compute the estimated GRE scores for GPA = 7.5 and 8.5..

Subject	11	12	13	14	15	16	17	18	19	20
X	8.3	8.6	9.2	9.8	8.0	7.8	9.4	9.0	7.2	8.6
y	2300	2250	2380	2400	2000	2100	2360	2350	2000	2260

8. A horse was subject to the test of how many minutes it takes to reach a point from the starting point. The horse was made to carry luggage of various weights on 10 trials.. The data collected are presented below in the table. Find the regression equation between the load and the time taken to reach the goal. Estimate the time taken for the loads of 35 Kgs , 23 Kgs, and 9 Kgs. Are the answers in agreement with your intuitive feelings? Justify.

Trial Number	1	2	3	4	5	6	8	8	9	10
Weight (in Kgs)	11	23	16	32	12	28	29	19	25	20
Time taken (in mins)	13	22	16	47	13	39	43	21	32	22

9. A study was conducted to find whether there is any relationship between the weight and blood pressure of an individual. The following set of data was arrived at from a clinical study.

Serial Number	1	2	3	4	5	6	8	8	9	10
Weight	78	86	72	822	80	86	84	89	68	71
Blood Pressure	140	160	134	144	180	176	174	178	128	132

10. It is assumed that achievement test scores should be correlated with student's classroom performance. One would expect that students who consistently perform well in the classroom (tests, quizzes, etc.) would also perform well on a standardized achievement test (0 - 100 with 100 indicating high achievement (x)). A teacher decides to examine this hypothesis. At the end of the academic year, she computes a correlation between the students achievement test scores (she purposefully did not look at this data until after she submitted students grades) and the overall g.p.a.(y) for each student computed over the entire year. The data for her class are provided below.

X	98	96	94	88	01	77	86	71	59	63	84	79	75	72	86	85	71	93	90	62
y	3.6	2.7	3.1	4.0	3.2	3.0	3.8	2.6	3.0	2.2	1.7	3.1	2.6	2.9	2.4	3.4	2.8	3.7	3.2	1.6

- Compute the correlation coefficient.
- What does this statistic mean concerning the relationship between achievement test performance and g.p.a.?
- What percent of the variability is accounted for by the relationship between the two variables and what does this statistic mean?
- What would be the slope and y-intercept for a regression line based on this data?
- If a student scored a 93 on the achievement test, what would be their predicted G.P.A.? If they scored a 74? A 88?

11. With the growth of internet service providers, a researcher decides to examine whether there is a correlation between cost of internet service per month (rounded to the nearest dollar) and degree of customer satisfaction (on a scale of 1 - 10 with a 1 being not at all satisfied and a 10 being extremely satisfied). The researcher only includes programs with comparable types of services. A sample of the data is provided below.

Cost of internet (in \$)	11	18	17	15	9	5	12	19	22	25
<u>satisfaction</u>	6	8	10	4	9	6	3	5	2	10

- Compute the correlation coefficient.
- What does this statistic mean concerning the relationship between amount of money spent per month on internet provider service and level of customer satisfaction?

- c) What percent of the variability is accounted for by the relationship between the two variables and what does this statistic mean?

12. It is hypothesized that there are fluctuations in norepinephrine (NE) levels which accompany fluctuations in affect with bipolar affective disorder (manic-depressive illness). Thus, during depressive states, NE levels drop; during manic states, NE levels increase. To test this relationship, researchers measured the level of NE by measuring the metabolite 3-methoxy-4-hydroxyphenylglycol (MHPG in micro gram per 24 hour) in the patient's urine experiencing varying levels of mania/depression. Increased levels of MHPG are correlated with increased metabolism (thus higher levels) of central nervous system NE. Levels of mania/depression were also recorded on a scale with a low score indicating increased mania and a high score increased depression. The data is provided below.

<u>MHPG</u>	980	1209	1403	1950	1814	1280	1073	1066	880	776
<u>Affect</u>	22	26	8	10	5	19	26	12	23	28

- a) Compute the correlation coefficient.
 b) What does this statistic mean concerning the relationship between MHPG levels and affect?
 c) What percent of the variability is accounted for by the relationship between the two variables?
 d) What would be the slope and y-intercept for a regression line based on this data?
 e) What would be the predicted affect score if the individual had an MHPG level of 1100? of 950? of 700?
13. The table below contains 25 cases -- the mother's weight in kilograms and the infant's birth weight in grams. Does this data suggest some relationship between the mother's weight and the infant's birth weight Why would such a relationship be important

. Prepregnancy Weights of Mothers and Birthweights of their Infants		
Case Number	Mother's Weight (kg)	Infant's Birthweight (g)
1	49.4	3515
2	63.5	3742
3	68.0	3629
4	52.2	2680
5	54.4	3006
6	70.3	4068
7	50.8	3373
8	73.9	4124
9	65.8	3572
10	54.4	3359
11	73.5	3230
12	59.0	3572
13	61.2	3062
14	52.2	3374
15	63.1	2722
16	65.8	3345
17	61.2	3714
18	55.8	2991

19	61.2	4026
20	56.7	2920
21	63.5	4152
22	59.0	2977
23	49.9	2764
24	65.8	2920
25	43.1	2693

5 PROBABILITY

5.1 What is Probability?

Probability theory is the branch of mathematics that studies the possible outcomes of given events together with the outcomes' relative likelihoods and distributions. In common usage, the word “probability” is used to mean the chance that a particular event (or set of events) will occur expressed on a linear scale from 0 (impossibility) to 1 (certainty). Factually, It is the study of random or indeterministic experiments eg tossing a coin or rolling a die. If we roll a die, we are certain it will come down but we are uncertain which face will show up. Ie the face showing up is indeterministic. Probability is a way of summarizing the uncertainty of statements or events. It gives a numerical measure for the degree of certainty (or degree of uncertainty) of the occurrence of an event.

We often use P to represent a probability Eg P(rain) would be the probability that it rains. In other cases Pr(.) is used instead of just P(.).

Definitions

- Experiment: A process by which an observation or measurement is obtained. Eg tossing a coin or rolling a die.
- Outcome: Possible result of a random experiment. Eg a 6 when a die is rolled once or a head when a coin is tossed.
- Sample space: Also called the probability space and it is a collection or set of all possible outcomes of a random experiment. Sample space is usually denoted by S or Ω or U
- Event: it's a subset of the sample space. Events are usually denoted by upper case letters.

Suppose the sample space S consists of $n(S)$ equally likely outcomes and $n(E)$ of those are favourable for an event E then probability of an event E is the ratio of the number of favourable outcomes $n(E)$ to the total number of all possible outcomes $n(S)$ ie

$$P(E) = \frac{\text{number of favourable outcomes}}{\text{total number of possible outcomes}} = \frac{n(E)}{n(S)}$$

5.2 Approaches to Probability

There are three ways to define probability, namely classical, empirical and subjective probability.

5.2.1 Classical probability

Classical or theoretical probability is used when each outcome in a sample space is equally likely to occur. The underlying idea behind this view of probability is symmetry. Ie if the sample space contains n outcomes that are fairly likely then $P(\text{one outcome}) = 1/n$.

The classical probability for an event A is given by

$$P(A) = \frac{\text{Number of outcomes in } A}{\text{Total number of outcomes in } S} = \frac{n(A)}{n(S)}$$

Eg Roll a die and observe that $P(A) = P(\text{rolling a 3}) = \frac{1}{6}$.

Example

A fair die, with faces numbered 1 to 6, is rolled once, write down the sample space S hence find the probability that the score showing up is ; a) a multiple of 3 b) a prime number.

Solution

$S = \{1, 2, 3, 4, 5, 6\}$ Multiples of 3 are 3 and 6 while prime numbers are 2, 3 and 5

Thus $P(\text{Multiple of 3}) = \frac{2}{6} = \frac{1}{3}$ and $P(\text{prime number}) = \frac{3}{6} = \frac{1}{2}$

5.2.2 Frequentist or Empirical probability

When the outcomes of an experiment are not equally likely, we can conduct experiments to give us some idea of how likely the different outcomes are. For example, suppose we are interested in measuring the probability of producing a defective item in a manufacturing process. The probability could be measured by monitoring the process over a reasonably long period of time and calculating the proportion of defective items.

In a nut shell Empirical (or frequentist or statistical) probability is based on observed data.

The empirical probability of an event A is the relative frequency of event A , that is

$$P(A) = \frac{\text{Frequency of event } A}{\text{Total number of observations}}$$

Example 1

The following are the counts of fish of each type, that you have caught before.

Fish Types	Blue gill	Red gill	Crappy	Total
No of times caught	13	17	10	40

Estimate the probability that the next fish you catch will be a Blue gill.

$$P(\text{Blue gill}) = \frac{13}{40} = 0.325$$

Example 2

A girl lists the number of male and female children her parent and her parent's brothers and sisters have. Her results were as tabulated below

	Males	Females			
Her parents	2	5	Her mother's brothers	4	8
Her mother's sisters	6	8	Her father's sisters	5	8
			Her father's brothers	7	7

	Totals	24	36
--	--------	----	----

- d) Find the probability that, if the girl has children of her own, the 1st born will be a girl.

- e) If the girl eventually has 10 children, how many are likely to be males?

Solution

- a) Following the family pattern, $P(\text{1st born will be a girl}) = \frac{36}{60} = 0.6$
b) 60% of the children will be females \Rightarrow 40% will be males. Thus 4 out of 10 children are likely to be males.

Remark: The empirical probability definition has a weakness that it depends on the results of a particular experiment. The next time this experiment is repeated, you are likely to get a somewhat different result. However, as an experiment is repeated many times, the empirical probability of an event, based on the combined results, approaches the theoretical probability of the event.

5.2.3 Subjective Probability:

Subjective probabilities result from intuition, educated guesses, and estimates. For example: given a patient's health and extent of injuries a doctor may feel that the patient has a 90% chance of a full recovery.. Subjectivity means two people can assign different probabilities to the same event.

Regardless of the way probabilities are defined, they always follow the same laws, which we will explore in the following Section.

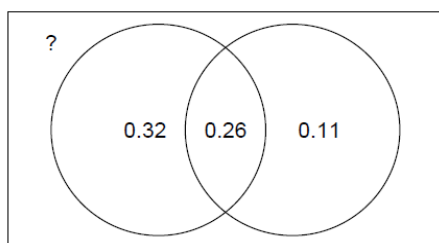
Exercise

- 1) What is the probability of getting a total of 7 or 11, when two dice are rolled?
- 2) Two cards are drawn from a pack, without replacement. What is the probability that both are greater than 2 and less than 8?
- 3) A permutation of the word "white" is chosen at random. Find the probability that it begins with a vowel. Also find the probability that it ends with a consonant.
- 4) Find the probability that a leap year will have 53 Sundays.
- 5) Two tetrahedral (4-sided) symmetrical dice are rolled, one after the other. Find the probability that;
 - a) both dice will land on the same number.
 - b) each die will land on a number less than 3.
 - c) the two numbers will differ by at most 1.

Will the answers change if we rolled the dice simultaneously?

Ways to represent probabilities:

- 1) *Venn diagram*; We may write the probabilities inside the elementary pieces within a Venn diagram. For example, $P(AB') = 0.32$ and $P(A) = P(AB) + P(AB') = 0.58$ [why?] The relative sizes of the pieces do not have to match the numbers.

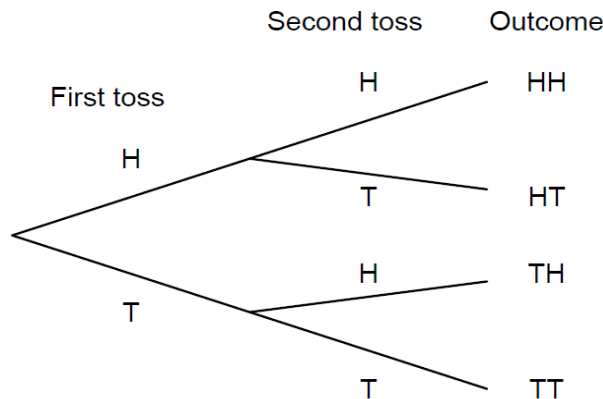


- 2) *Two-way table*; This is a popular way to represent statistical data. The cells of the table correspond to the intersections of row and column events. Note that the contents of the table

add up accross rows and columns of the table. The bottom-right corner of the table contains $P(S) = 1$

	B	B'	
A	0.26	0.32	0.58
A'	0.11	?	0.42
	0.37	0.63	1

Tree diagram; Tree diagrams or probability trees are simpler clear ways of representing probabilistic information. A tree diagram may be used to show the sequence of choices that lead to the complete description of outcomes. For example, when tossing two coins, we may represent this as follows



A tree diagram is also often useful for representing conditional probabilities

5.3 Review of set notation

Complement: The complement of event A, (denoted A'), is the set of all outcomes in a sample that are not included in the event A.

Intersection of events: The event $A \cap B$ (or simply AB) read as 'A intersection B' consists of outcomes that are contained within both events A and B. The probability of this event is the probability that both events A and B occur [but not necessarily at the same time]. Here after we will abbreviate intersection as AB .

Unions of Events: The event $A \cup B$ read as 'A union B' consists of the outcomes that are contained within at least one of the events A and B. The probability of this event $P(A \cup B)$; is the probability that events A and/or B occurs.

Set notation

Suppose a set S consists of points labelled 1, 2, 3 and 4. We denote this by $S = \{1, 2, 3, 4\}$. . If

$A = \{1, 2\}$ and $B = \{2, 3, 4\}$, then A and B are subsets of S, denoted by

$A \subset S$ and $B \subset S$ (B is contained in S). We denote the fact that 2 is an element of A by $2 \in A$.

The union of A and B, $A \cup B = \{1, 2, 3, 4\}$. If $C = \{4\}$, then $A \cup C = \{1, 2, 4\}$. The intersection $A \cap B = AB = \{2\}$: The complement of A, is $A' = \{3, 4\}$.

Distributive laws; $A \cap (B \cup C) = AB \cup AC$ and $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$

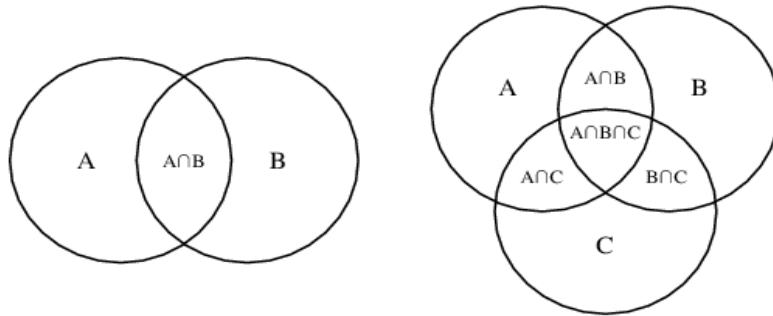
De Morgan's Law; $(A \cup B)' = A'B'$ and $(AB)' = A' \cup B'$

Venn diagram

Venn diagram is a diagram that shows all possible logical relations between a finite collection of sets.

Venn diagram is often used to illustrate the relations between sets (events).

The sets A and B are represented as circles; operations between them (intersections, unions and complements) can also be represented as parts of the Venn diagram. The entire sample space S is the bounding box. See Figure 2.1



Exercise

- Use the Venn diagrams to illustrate Distributive laws and De Morgan's law.
- Simplify the following (Draw the Venn diagrams to visualize)
 - $(A')'$
 - $(AB)' \cup A$
 - $AB \cup AB'$
 - $(A \cup B) \cap B$
- Represent by set notation and exhibit on a Venn diagram the following events
 - both A and B occur
 - exactly one of A, B occurs
 - A and B but not C occurs
 - at least one of A, B, C occurs
 - at most one of A, B, C occurs
- The sample space consists of eight capital letters (outcomes), A, B, C, ..., H. Let V be the event that the letter represents a vowel, and L be the event that the letter is made of straight lines. Describe the outcomes that comprise
 - VL
 - $V \cup L'$
 - $V'L'$
- Out of all items sent for refurbishing, 40% had mechanical defects, 50% had electrical defects, and 25% had both. Denoting A = fan item has a mechanical defect and B = fan item has an electrical defect, fill the probabilities into the Venn diagram and determine the quantities listed below.
 - $P(A)$
 - $P(AB)$
 - $P(A'B)$
 - $P(A'B')$
 - $P(A \cup B)$
 - $P(A \cup B')$
 - $P[(A \cup B)']$
- A sample of mutual funds was classified according to whether a fund was up or down last year (A and A') and whether it was investing in international stocks (B and B'). The probabilities of these events and their intersections are represented in the two-way table below. Fill in all the question marks hence find the probability of $A \cup B$

	B	B'	
A	0.33	?	?
A'	?	?	0.52
	0.64	?	1

5.4 Rules of Probability

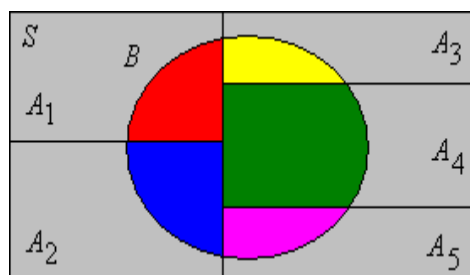
- For an experiment with a sample space $S = \{E_1, E_2, \dots, E_n\}$ we can assign probabilities

$$P(E_1), P(E_2), \dots, P(E_n) \text{ provided that } 0 \leq P(E_i) \leq 1 \text{ and } P(S) = \sum_{i=1}^n P(E_i) = 1$$

Remark:

- a) If a set (event) A consists of outcomes E_1, E_2, \dots, E_k , then $P(A) = \sum_{i=1}^k P(E_i)$
- b) If $E = S$ then $P(E) = P(S) = 1$ and If E has no elements, (ie if E is empty)
 $\Rightarrow E = \phi$ or $\{ \}$, then $P(E) = P(\phi) = 0$
- 2) If E is an event in the sample space S , then E' (called the complement of E) is an event in S but outside E . $P(E) + P(E') = 1 \Rightarrow P(E') = 1 - P(E)$
- 3) If the sample space S contains n disjoint events E_1, E_2, \dots, E_n , then
- $$P(E_1) + P(E_2) + \dots + P(E_n) = \sum_{i=1}^n P(E_i) = 1$$
- 4) Let A and B be two events such that $A \subseteq B$, then $P(A) \leq P(B)$
- 5) For any two events A and B , $P(A \cup B) = P(A) + P(B) - P(AB)$ where $P(AB) = P(A \cap B)$.
 Extension of this rule leads to the **Inclusion-Exclusion Principle**. This principle is a way to extend the general addition rule to 3 or more events. Here we will limit it to 3 events.
 $P(A \cup B \cup C) = P(A) + P(B) + P(C) - P(AB) - P(AC) - P(BC) + P(ABC)$
- 6) **Law of Partitions:** The law of partitions is a way to calculate the probability of an event. Let A_1, A_2, \dots, A_k form a partition of the sample space Ω . then, for any events B ,

$$P(B) = P(A_1 B) + P(A_2 B) + \dots + P(A_k B) = \sum_{i=1}^k P(A_i B)$$



Example 1

The Probability that John passes a Maths exam is $\frac{4}{5}$ and that he passes a Chemistry exam is $\frac{5}{6}$. If the probability that he passes both exams is $\frac{3}{4}$, find the probability that he will pass at least one exam.

Solution

Let M be the event that John passes Math exam, and C be the event that John passes Chemistry exam.

$$P(\text{John passes at least one exam}) = P(M \cup C) = P(M) + P(C) - P(MC) = \frac{4}{5} + \frac{5}{6} - \frac{3}{4} = \frac{53}{60}$$

Example 2

A fair die, with faces numbered 1 to 6, is rolled twice and the sum of the scores showing up noted. Let A be the event that the sum of the scores is greater than 7, B be the event that the sum of the scores is a multiple of 3 and C be the event that the sum of the scores is a prime number. Show that $P(A \cup B) = P(A) + P(B) - P(AB)$ and also find $P(A \cup C)$, $P(BC)$ and $P(BC')$

Solution

+	1	2	3	4	5	6
1	2	3	4	5	6	7
2	3	4	5	6	7	8

3	4	5	6	7	8	9
4	5	6	7	8	9	10
5	6	7	8	9	10	11
6	7	8	9	10	11	12

1. $P(A) = \frac{15}{36} = \frac{5}{12}$ and $P(B) = \frac{12}{36} = \frac{1}{3}$
 $A \cap B$ means the set of all multiples of 3 that are greater than 7. Clearly
 $P(A \cap B) = \frac{5}{36}$

$A \cup B$ means the set of all values that are multiples of 3 and/or greater than 7.

Clearly

$$P(A \cup B) = \frac{22}{36} = \frac{11}{18} = P(A) + P(B) - P(AB)$$

2. $P(C) = \frac{15}{36} = \frac{5}{12}$

$A \cap C$ means the set of all multiples of 3 that are prime number. Clearly $P(A \cap C) = \frac{2}{36} = \frac{1}{18}$

$$P(A \cup C) = P(A) + P(C) - P(AC) = \frac{5}{12} + \frac{5}{12} - \frac{1}{18} = \frac{7}{9}$$

$B \cap C$ means the set of all greater than 7 that are prime numbers. Clearly $P(BC) = \frac{2}{36} = \frac{1}{18}$

$B \cap C'$ means the set of all greater than 7 that are not prime numbers. Clearly $P(BC') = \frac{11}{36}$

Exercise

- Which of the following is a probability function defined on $S = \{E_1, E_2, E_3\}$
 - $P(E_1) = \frac{1}{4}, P(E_2) = \frac{1}{3}$ and $P(E_3) = \frac{1}{2}$
 - $P(E_1) = \frac{1}{3}, P(E_2) = \frac{1}{6}$ and $P(E_3) = \frac{1}{2}$
 - $P(E_1) = \frac{2}{3}, P(E_2) = -\frac{1}{3}$ and $P(E_3) = \frac{2}{3}$
 - $P(E_1) = 0, P(E_2) = \frac{1}{3}$ and $P(E_3) = \frac{2}{3}$
- As a foreign language, 40% of the students took Spanish and 30% took French, while 60% took at least one of these languages. What percent of students took both Spanish and French?
- In a class of 100 students, 30 are in mathematics. Moreover, of the 40 females in the class, 10 are in Mathematics. If a student is selected at random from the class, what is the probability that the student will be a male or be in mathematics?
- The probability that a car stopped at a road brook will have faulty breaks is 0.23, the probability that it will have badly worn out tyres is 0.24 and the probability that it will have faulty breaks and/or badly worn out tyres is 0.38. Find the probability that a car which has just been stopped will have both faulty breaks and badly worn out tyres.
- Given two events A and B in the same sample space such that $P(A) = 0.59$, $P(B) = 0.3$ and $P(AB) = 0.21$. Find; a) $P(A \cup B)$ b) $P(A'B)$ c) $P(AB')$ d) $P(A \cup B')$
- Let A and B be two events in the same sample space such that $P(A \cup B) = \frac{3}{4}$, $P(B') = \frac{2}{3}$ and $P(AB) = \frac{1}{4}$. Find $P(B)$, $P(A)$ and $P(AB')$
- Suppose that $P(A) = 0.4$, $P(B) = 0.5$ and $P(A \cap B) = 0.2$ Find; a) $P(A \cup B)$ b) $P(A'B')$ c) $P[A' \cap (A \cup B)]$ d) $P[A \cup (A'B)]$
- A die is loaded such that even numbers are twice as likely as odd numbers. Find the probability that for a single toss of this die the spot showing up is greater than 3
- A point is selected at random inside an equilateral triangle of sides 3 units. Find the probability that its distance to any corner is greater than 1 unit.

Definition: (Odds of an event)

It's the ratio of the probability of an event occurring to that of the event not happening. If A is an event then the odds of A is given by $\frac{P(A)}{1-P(A)} = \frac{P(A)}{1-P(A)}$

Example

Find $P(A)$ and $P(A')$ if the odds of event A is $\frac{5}{4}$

Solution

$$\frac{P(A)}{1-P(A)} = \frac{5}{4} \Rightarrow 4(1-P(A)) = 5P(A) \Rightarrow 4 = 9P(A) \Rightarrow P(A) = \frac{4}{9} \text{ and } P(A') = \frac{5}{9}$$

Question: Find $P(E)$ and $P(E')$ if the odds of event E is (i) $\frac{3}{4}$ (ii) $\frac{a}{b}$

5.5 Relationship between Events

- Compound event: Two or more events combined together. Eg AB is a compound event
- Mutually exclusive events: Two events A and B are said to be mutually exclusive if they cannot occur simultaneously. That is if the occurrence of A totally excludes the occurrence of B. Effectively events A and B are said to be mutually exclusive if they disjoint. ie $A \cap B = \emptyset \Rightarrow P(AB) = 0$
- Exhaustive events: Events whose union equals the sample space.
- Independent events: Two events A and B are said to be independent if the occurrence of A does not affect the occurrence of B. If events A and B are independent then $P(AB) = P(A) \times P(B)$

Remark: Three events A, B and C are said to be jointly independent if and only if

- (i) $P(AB) = P(A) \times P(B)$, $P(AC) = P(A) \times P(C)$ and $P(BC) = P(B) \times P(C)$ (ie they are pairwise independent) and
(ii) if $P(ABC) = P(A) \times P(B) \times P(C)$

Note it does not necessarily mean that if events A, B and C are pairwise independent then they are jointly independent

Example 1

Roll a fair die twice and define A to be the event that the sum of the scores showing up is greater than 7, B be the event that the sum of the scores showing up is a multiple of 3 and C be the event that the sum of the scores showing up is a prime number. Which of the events A, B and C are independent? Are the 3 events jointly independent?

Solution

From the above example, $P(A) = P(C) = \frac{5}{12}$, $P(B) = \frac{1}{3}$, $P(AB) = \frac{5}{36}$ and $P(AC) = P(BC) = \frac{1}{18}$
Since $P(AB) \neq 0$ events A and B are not mutually exclusive. Similarly events A & C and B and C are not mutually exclusive

$$P(A) \times P(B) = \frac{5}{12} \times \frac{1}{3} = \frac{5}{36} = P(AB) \Rightarrow A \text{ and } B \text{ are independent events.}$$

$$P(A) \times P(C) = \frac{5}{12} \times \frac{5}{12} = \frac{25}{144} \neq P(AC) \Rightarrow A \text{ and } C \text{ are dependent events.}$$

$$P(B) \times P(C) = \frac{1}{3} \times \frac{5}{12} = \frac{5}{36} \neq P(BC) \Rightarrow B \text{ and } C \text{ are dependent events.}$$

The 3 events are not jointly independent since pairwise independence is not satisfied.

Example 2

Three different machines in a factory have the following probabilities of breaking down during a shift.

Machine	A	B	C
probability	$\frac{4}{15}$	$\frac{3}{10}$	$\frac{2}{11}$

Find the probability that in a particular shift,;

- a) All the machines will break down
b) None of the machines will break down.

Solution

Since the events of breaking down of machines are independent, probability that all the machines will break down is given by $P(ABC) = \frac{4}{15} \times \frac{3}{10} \times \frac{2}{11} = \frac{4}{275}$

The probability that none of the machines will break down is given by

$$P(\overline{ABC}) = \frac{11}{15} \times \frac{7}{10} \times \frac{9}{11} = \frac{21}{50}$$

Exercises

- 1) In a game of archery the probability that A hits the target is $\frac{1}{3}$ and the probability that B hits the target is $\frac{2}{5}$. What is the probability that the target will be hit?
- 2) Toss a fair coin 3 times and let A be the event that two or more heads appears, B be the event that all outcomes are the same and C be the event that at most two tails appears. Which of the events A, B and C are independent? Are the 3 events jointly independent?
- 3) A fair coin and a fair die are rolled together once. Let A be the event that a head and an even number appears, B be the event that a prime number appears and C be the event that a tail and an odd number appears.
 - a) Express explicitly the event that i) A and B occurs ii) Only B occurs iii) B and C occur
 - b) Which of the events A, B and C are independent and which ones are mutually exclusive?
- 4) A die is loaded so that the probability of a face showing up is proportional to the face number. Write down the probability of each sample point. If A is the event that an even number appears, B is the event that a prime number appears and C is the event that an odd number appears.
 - a) Find the probability that: i) A and/or B occurs ii) A but not B occurs iii) B and C occurs d) A and/or C occurs
 - b) Which of the events A, B and C are independent and which ones are mutually exclusive?

Theorem 1: If events A and B are independent, then A and B' are also independent

Proof

Decomposing A into two disjoint events AB and AB' . We can write

$P(A) = P(AB) + P(AB') \Rightarrow P(AB') = P(A) - P(AB) = P(A) - P(A) \times P(B)$ since events A and B are independent. Thus $P(AB') = P(A)[1 - P(B)] = P(A) \times P(B') \Rightarrow A$ and B' are independent

Theorem 2: If events A and B are independent, then A' and B' are also independent

Proof

Decomposing B' into two disjoint events AB' and $A'B'$. We can write

$P(B') = P(AB') + P(A'B') \Rightarrow P(A'B') = P(B') - P(AB') = P(B') - P(A) \times P(B')$ since events A and B' are independent. (From theorem 1 above) Thus

$P(A'B') = [1 - P(A)]P(B') = P(A') \times P(B') \Rightarrow A'$ and B' are also independent

5.6 Counting Rules useful in Probability

In some experiments it is helpful to list the elements of the sample space systematically by means of a tree diagram,. In many cases, we shall be able to solve a probability problem by counting the number of points in the sample space without actually listing each element.

Theorem (Multiplication principle)

If one operation can be performed in n_1 ways, and if for each of these a second operation can be performed in n_2 ways, then the two operations can be performed together in $n_1 n_2$ ways.

Eg How large is the sample space when a pair of dice is thrown?

Solution; The first die can be thrown in $n_1 = 6$ ways and the second in

$n_2 = 6$ Ways. Therefore, the pair of dice can land in $n_1 n_2 = 36$ possible ways.

The above theorem can naturally be extended to more than two operations: if we have n_1, n_2, \dots, n_k consequent choices, then the total number of ways is $n_1 \times n_2 \times \dots \times n_k$

Permutations

Permutations refer to an arrangement of objects when the order matters (for example, letters in a word). The number of permutations of n distinct objects taken r at a time is ${}_nP_r = \frac{n!}{(n-r)!}$

Example

From among ten employees, three are to be selected to travel to three out-of-town plants A, B, and C, one to each plant. Since the plants are located in different cities, the order in which the employees are assigned to the plants is an important consideration. In how many ways can the assignments be made?

Solution;

Because order is important, the number of possible distinct assignments is ${}_{10}P_3 = 720$

In other words, there are ten choices for plant A, but then only nine for plant B, and eight for plant C. This gives a total of $10(9)(8)$ ways of assigning employees to the plants.

Combinations

The term combination refers to the arrangement of objects when order does not matter. For example, choosing 4 books to buy at the store in any order will leave you with the same set of books. The number of distinct subsets or combinations of size r that can be selected from n

distinct objects, $(r \leq n)$, is given by ${}_nC_r = \frac{n!}{r!(n-r)!}$

Example 1

In the previous example, suppose that three employees are to be selected from among the ten available to go to the same plant. In how many ways can this selection be made?

Solution

Here, order is not important; we want to know how many subsets of size $r = 3$ can be selected from $n = 10$ people. The result is ${}_{10}C_3 = 120$

Example 2

In a poker hand consisting of 5 cards, find the probability of holding 2 aces and 3 jacks.

Solution

The number of ways of being dealt 2 aces from 4 is ${}_4C_2 = 6$ and the number of ways of being dealt 3 jacks from 4 is ${}_4C_3 = 4$

The total number of 5-card poker hands, all of which are equally likely is ${}_{52}C_5 = 2,598,960$

Hence, the probability of getting 2 aces and 3 jacks in a 5-card poker hand is $P(C) =$

$$P(C) = \frac{6 \times 4}{2,598,960}$$

Example 3

A university warehouse has received a shipment of 25 printers, of which 10 are laser printers and 15 are inkjet models. If 6 of these 25 are selected at random to be checked by a particular technician, what is the probability that; a) exactly 3 of these selected are laser printers? b) at least 3 inkjet printers?

Solution

First choose 3 of the 15 inkjet and then 3 of the 10 laser printers.

There are ${}_{15}C_3$ and ${}_{10}C_3$ ways to do it, and therefore

$$P(\text{exactly 3 of the 6}) = \frac{{}^{15}C_3 \times {}^{10}C_3}{{}^{25}C_6} = 0.3083$$

$$P(\text{at least 3}) = \frac{{}^{15}C_3 \times {}^{10}C_3}{{}^{25}C_6} + \frac{{}^{15}C_4 \times {}^{10}C_2}{{}^{25}C_6} + \frac{{}^{15}C_5 \times {}^{10}C_1}{{}^{25}C_6} + \frac{{}^{15}C_6 \times {}^{10}C_0}{{}^{25}C_6} = 0.8530$$

Exercises

- 1) An incoming lot of silicon wafers is to be inspected for defectives by an engineer in a microchip manufacturing plant. Suppose that, in a tray containing 20 wafers, 4 are defective. Two wafers are to be selected randomly for inspection. Find the probability that neither is defective.
- 2) A person draws 5 cards from a shuffled pack of cards. Find the probability that the person has at least 3 aces. Find the probability that the person has at least 4 cards of the same suit.
- 3) A California licence plate consists of a sequence of seven symbols: number, letter, letter, letter, number, number, number, where a letter is any one of 26 letters and a number is one among 0, 1, ... 9. Assume that all licence plates are equally likely. What is the probability that;
 - a) all symbols are different?
 - b) all symbols are different and the first number is the largest among the numbers?
- 4) A bag contains 80 balls numbered 1.... 80. Before the game starts, you choose 10 different numbers from amongst 1.... 80 and write them on a piece of paper. Then 20 balls are selected (without replacement) out of the bag at random. What is the probability that;
 - a) all your numbers are selected?
 - b) none of your numbers is selected?
 - c) exactly 4 of your numbers are selected?
- 5) A full deck of 52 cards contains 13 hearts. Pick 8 cards from the deck at random (a) without replacement and (b) with replacement. In each case compute the probability that you get no hearts.
- 6) Three people enter the elevator on the basement level. The building has 7 floors. Find the probability that all three get off at different floors.
- 7) In a group of 7 people, each person shakes hands with every other person. How many handshakes did occur?
- 8) A marketing director considers that there's "overwhelming agreement" in a 5-member focus group when either 4 or 5 people like or dislike the product. If, in fact, the product's popularity is 50% (so that all outcomes are equally likely), what is the probability that the focus group will be in "overwhelming agreement" about it? Is the marketing director making a judgement error in declaring such agreement "overwhelming"?
- 9) A die is tossed 5 times. Find the probability that we will have 4 of a kind.
- 10) A tennis tournament has $2n$ participants, n Swedes and n Norwegians. First, n people are chosen at random from the $2n$ (with no regard to nationality) and then paired randomly with the other n people. Each pair proceeds to play one match. An outcome is a *set* of n (ordered) pairs, giving the winner and the loser in each of the n matches. (a) Determine the number of outcomes. (b) What do you need to assume to conclude that all outcomes are equally likely? (c) Under this assumption, compute the probability that all Swedes are the winners.
- 11) A group of 18 Scandinavians consists of 5 Norwegians, 6 Swedes, and 7 Finns. They are seated at random around a table. Compute the following probabilities: (a) that all the Norwegians sit together, (b) that all the Norwegians and all the Swedes sit together, and (c) that all the Norwegians, all the Swedes, and all the Finns sit together.

- 12) In a lottery, 6 numbers are drawn out of 45. You hit a jackpot if you guess all 6 numbers correctly, and get \$400 if you guess 5 numbers out of 6. What are the probabilities of each of those events?
- 13) There are 21 Bachelor of Science programs at New Mexico Tech. Given 21 areas from which to choose, in how many ways can a student select:
 - a) A major area and a minor area?
 - b) A major area and first and second minor?
- 14) From a box containing 5 chocolates and 4 hard candies, a child takes a handful of 4 (at random). What is the probability that exactly 3 of the 4 are chocolates?
- 15) If a group consist of 8 men and 6 women, in how many ways can a committee of 5 be selected if:
 - a) The committee is to consist of 3 men and 3 women.
 - b) There are no restrictions on the number of men and women on the committee.
 - c) There must at least one man.
 - d) There must be at least one of each sex.
- 16) Suppose we have a lot of 40 transistors of which 8 are defective. If we sample without replacement, what is the probability that we get 4 good transistors in the first 5 draws?
- 17) A housewife is asked to rank four brands A, B, C, and D of household cleaner according to her preference, number one being the one she prefers most, etc. she really has no preference among the four brands. Hence, any ordering is equally likely to occur.
 - a) Find the probability that brand A is ranked number one.
 - b) Find the probability that brand C is number one D is number 2 in the rankings.
 - c) Find the probability that brand A is ranked number one or number 2.
- 18) How many ways can one arrange the letters of the word ADVANTAGE so
- 19) that the three As are adjacent to each other?
- 20) Eight tires of different brands are ranked 1 to 8 (best to worst) according to mileage performance. If four of these tires are chosen at random by a customer, find the probability that the best tire among the four selected by the customer is actually ranked third among the original eight.

5.7 Conditional Probability and Independence

Humans often have to act based on incomplete information. If your boss has looked at you gloomily, you might conclude that something's wrong with your job performance. However, if you know that she just suffered some losses in the stock market, this extra information may change your assessment of the situation. Conditional probability is a tool for dealing with additional information like this.

Conditional probability is the probability of an event occurring given the knowledge that another event has occurred. The conditional probability of event A occurring, given that event B has occurred is denoted by $P(A/B)$ and is read "probability of A given B" and is given by

$$P(A/B) = \frac{P(AB)}{P(B)} \quad \text{provided } P(B) > 0 \quad \text{Similarly } P(B/A) = \frac{P(AB)}{P(A)} \quad \text{provided } P(A) > 0$$

$$\Rightarrow P(AB) = P(A/B) \times P(B) = P(B/A) \times P(A)$$

Remark: Another way to express independence is to say that the knowledge of B occurring does not change our assessment of $P(A)$. This means that if A and B are independent then $P(A/B) = P(A)$ and $P(B/A) = P(B)$

Example

In a large metropolitan area, the probability of a family owning a colour T.V , a computer or both 0.86, 0.35 and 0.29 respectively. What is the probability that a family chosen at random during a survey will own a colour T.V and/or a computer? Given that the family chosen at

random during a survey owns a colour T.V, what is the probability that it will own a computer?

Solution

Let T and C be the event of owning a colour T.V and a computer respectively. Then

$$P(T \cup C) = P(T) + P(C) - P(TC) = 0.86 + 0.35 - 0.29 = 0.92$$

$$P(C/T) - \frac{P(TC)}{P(T)} = \frac{0.29}{0.86} \approx 0.$$

Reduced sample space approach

In case when all the outcomes are equally likely, it is sometimes easier to find conditional probabilities directly, without having to apply the above equation. If we already know that B has happened, we need only to consider outcomes in B, thus reducing our sample space to B.

$$\text{Then, } P(A/B) = \frac{\text{Number of outcomes in } AB}{\text{Number of outcomes in } B}$$

For example, $P(\text{a die is } 3 / \text{a die is odd}) = \frac{1}{3}$ and $P(\text{a die is } 4 / \text{a die is odd}) = 0$

Example

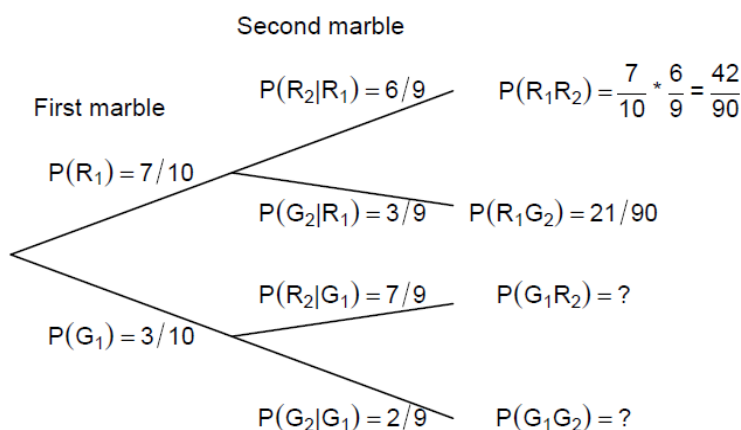
Let $A = \{\text{a family has two boys}\}$ and $B = \{\text{a family of two has at least one boy}\}$ Find $P(A/B)$

Solution

The event B contains the following outcomes: $B = \{(B, B), (B, G), (G, B)\}$ and. Only one of these is in A. Thus, $P(A/B) = \frac{1}{3}$. However, if I know that the family has two children, and I see one of the children and it's a boy, then the probability suddenly changes to $1/2$. There is a subtle difference in the language and this changes the conditional probability

5.7.1 Tree Diagrams in conditional probability

Suppose we are drawing marbles from a bag that initially contains 7 red and 3 green marbles. The drawing is without replacement that is after we draw the first marble, we do not put it back. Let's denote the events $R_1 = \{\text{the first marble is red}\}$ $R_2 = \{\text{the second marble is red}\}$ $G_1 = \{\text{the first marble is green}\}$ and so on. Let's fill out the tree representing the consecutive choices.



The conditional probability $P(R_2 / R_1)$ can be obtained directly from reasoning that after we took the first red marble, there remain 6 red and 3 green marbles. On the other hand, we could

use the formula to get $P(R_2/R_1) = \frac{P(R_1 R_2)}{P(R_1)} = \frac{\frac{42}{90}}{\frac{7}{10}} = \frac{2}{3}$ where the probability $P(R_2/R_1)$

{same as $P(R_1 R_2)$ } can be obtained from counting the outcomes $P(R_1 R_2) = \frac{{}^7C_2}{{}^{10}C_2} = \frac{7}{15}$

Question: Find $P(R_2)$ and $P(R_1/R_2)$.

Example 2

Suppose that of all individuals buying a certain digital camera, 60% include an optional memory card in their purchase, 40% include a set of batteries, and 30% include both a card and batteries. Consider randomly selecting a buyer and let $A = \{\text{memory card purchased}\}$ and $B = \{\text{battery purchased}\}$. Then find $P(A/B)$ and $P(B/A)$.

Solution

From given information, we have $P(A) = 0.60$, $P(B) = 0.40$ and

$P(\text{both purchased}) = P(A \cap B) = 0.30$

Given that the selected individual purchased an extra battery, the probability that an optional

card was also purchased is $P(A/B) = \frac{P(A \cap B)}{P(B)} = \frac{0.30}{0.40} = 0.75$

That is, of all those purchasing an extra battery, 75% purchased an optional memory card.

Similarly $P(\text{battery } | \text{ memory card}) = P(B/A) = \frac{P(A \cap B)}{P(A)} = \frac{0.30}{0.60} = 0.5$

Notice that $P(A/B) \neq P(A)$ and $P(B/A) \neq P(B)$, that is, the events A and B are dependent.

Remark: The tree diagram may become tedious especially when the tree grows beyond 4 stages. In such a case we can make use of binomial formula which is applicable when:

- The experiment's outcome can be classified into 2 categories success and failure with probabilities p and $1-p$ respectively
- The experiment is to be repeated n independent times
- Our interest is the number of successes

The probability of observing x successes out of n trials is given by:-

$$P(x) = {}_n C_x \times p^x (1-p)^{n-x} \text{ for } x = 0, 1, 2, \dots, n$$

Example

A fair coin is tossed 10 times. What is the probability of observing exactly 8 heads?

Solution

$n = 10$ $p = 0.5$ and $x = 8$ successes Therefore

$$P(X = 8) = {}_{10}C_8 \times 0.5^8 (1-0.5)^{10-8} = {}_{10}C_8 \times 0.5^{10} \approx 0.044$$

Exercises

- A pair of fair dice is rolled once. If the sum of the scores showing up is 6, find the probability that one of the dice shows a 2.
- A consumer research organisation has studied the services and warranty provided by 50 new car dealers in a certain city. Its findings are as follows

In Business for	Good services and a warranty	Poor services and a warranty
At least 10 years	16	4
Less than 10 years	10	20

If a person randomly selects one of these new car dealers ;

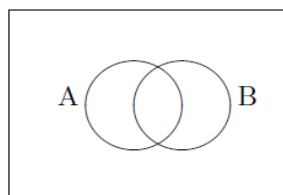
- a) What is the probability that he gets one who provides good services and a warranty
 - b) Who has been in business for at least 10 years, what is the probability that hehe provides good services and a warranty
 - c) What is the probability that one of these new car dealers who has been in business for less than 10 years will provide good services and a warranty?
- 3) Three machines A, B and C produces 50%, 30% and 20% respectively of the total number of items in a factory. The percentage of defective outputs of these machines are 3%, 4% and 5% respectively. If an item is selected at random:-
- a) Find the probability that it is defective
 - b) And found to be defective, what is the probability that it was produced by machine A?
- 4) A year has 53 Sundays. What is the conditional probability that it is a leap year?
- 5) The probability that a majority of the stockholders of a company will attend a special meeting is 0.5. If the majority attends, then the probability that an important merger will be approved is 0.9. What is the probability that a majority will attend and the merger will be approved?
- 6) Let events A, B have positive probabilities. Show that, if $P(A/B) = P(A)$ then also $P(B/A) = P(B)$.
- 7) The cards numbered 1 through 10 are placed in a hat, mixed up, and then one of the cards is drawn. If we are told that the number on the drawn card is at least five, what is the probability that it is ten?
- 8) In the roll of a fair die, consider the events $A = \{2, 4, 6\} = \text{"even numbers"}$ and $B = \{4, 5, 6\} = \text{"high scores"}$. Find the probability that die showing an even number given that it is a high score.
- 9) There are two urns. In the first urn there are 3 white and 2 black balls and in the second urn there 1 white and 4 black balls. From a randomly chosen urn, one ball is drawn. What is the probability that the ball is white?
- 10) The level of college attainment of US population by racial and ethnic group in 1998 is given in the following table

Racial or Ethnic Group	No of Adults (Millions)	%age with Associate's Degree	%age With Bachelor's Degree	%age with Graduate or Professional Degree
Native Americans	1.1	6.4	6.1	3.3
Blacks	16.8	5.3	7.5	3.8
Asians	4.3	7.7	22.7	13.9
Hispanics	11.2	4.8	5.9	3.3
Whites	132.0	6.3	13.9	7.7

The percentages given in the right three columns are conditional percentages.

- a) How many Asians have had a graduate or professional degree in 1998?
 - b) What percent of all adult Americans has had a Bachelor's degree?
 - c) Given that the person had an Associate's degree, what is the probability that the person was Hispanic?
- 11) The dealer's lot contains 40 cars arranged in 5 rows and 8 columns. We pick one car at random. Are the events $A = \{\text{the car comes from an odd-numbered row}\}$ and $B = \{\text{the car comes from one of the last 4 columns}\}$ independent? Prove your point of view.
- 12) You have sent applications to two colleges. If you are considering your chances to be accepted to either college as 60%, and believe the results are statistically independent, what is the probability that you'll be accepted to at least one? How will your answer change if you applied to 5 colleges?
- 13) In a high school class, 50% of the students took Spanish, 25% took French and 30% of the students took neither. Let A be the event that a randomly chosen student took Spanish,

and B be the event that a student took French. Fill in either the Venn diagram or a 2-way table and answer the questions:



	B	B'
A		
A'		

- Describe in words the meaning of the event AB' . Find the probability of this event.
 - Are the events A, B independent? Explain with numbers why or why not.
 - If it is known that the student took Spanish, what are the chances that she also took French?
- 14) One half of all female physicists are married. Among those married, 50% are married to other physicists, 29% to scientists other than physicists and 21% to non-scientists'. Among male physicists, 74% are married. Among them, 7% are married to other physicists, 11% to scientists other than physicists and 82% to non-scientists. What percent of all physicists are female? [Hint: This problem can be solved as is, but if you want to, assume that physicists comprise 1% of all population.]
- 15) Error-correcting codes are designed to withstand errors in data being sent over communication lines. Suppose we are sending a binary signal (consisting of a sequence of 0's and 1's), and during transmission, any bit may get flipped with probability p , independently of any other bit. However, we might choose to repeat each bit 3 times. For example, if we want to send a sequence 010, we will code it as 000111000. If one of the three bits flips, say, the receiver gets the sequence 001111000, he will still be able to decode it as 010 by majority voting. That is, reading the first three bits, 001, he will interpret it as an attempt to send 000. However, if two of the three bits are flipped, for example 011, this will be interpreted as an attempt to send 111, and thus decoded incorrectly. What is the probability of a bit being decoded incorrectly under this scheme?

5.8 Bayes' Rule

Events B_1, B_2, \dots, B_K are said to be a partition of the sample space S if the following two conditions are satisfied. i) $B_i B_j = \emptyset$ for each pair i, j and ii) $B_1 \cup B_2 \cup \dots \cup B_K = S$

This situation often arises when the statistics are available in subgroups of a population. For example, an insurance company might know accident rates for each age group B_i . This will give the company conditional probabilities $P(A/B_i)$ (if we denote $A = \{\text{event of accident}\}$).

Question: if we know all the conditional probabilities $P(A/B_i)$, how do we find the unconditional $P(A)$?

Consider a case when $k = 2$:

The event A can be written as the union of mutually exclusive events AB_1 and AB_2 , that is $A = AB_1 \cup AB_2$ it follows that $P(A) = P(AB_1) + P(AB_2)$

If the conditional probabilities $P(A/B_1)$ and $P(A/B_2)$ are known, that is

$$P(A/B_1) = \frac{P(AB_1)}{P(B_1)} \text{ and } P(A/B_2) = \frac{P(AB_2)}{P(B_2)} \text{ then } P(A) = P(B_1) \times P(A/B_1) + P(B_2) \times P(A/B_2)$$

Suppose we want to find probability of the form $P(B_i/A)$, which can be written as

$$P(B_i/A) = \frac{P(AB_i)}{P(A)} = \frac{P(B_i) \times P(A/B_i)}{P(A)} = \frac{P(B_i) \times P(A/B_i)}{P(B_1) \times P(A/B_1) + P(B_2) \times P(A/B_2)}$$

This calculation generalizes to $k > 2$ events as follows.

Theorem

If B_1, B_2, \dots, B_K form a partition of the sample space S such that $P(B_i) \neq 0$ for $i = 1, 2, \dots, k$; then for any event $A \subseteq S$,

$$P(A) = \sum_{i=1}^k P(AB_i) = \sum_{i=1}^k P(B_i) \times P(A/B_i) \quad \text{Subsequently, } P(B_i/A) = \frac{P(AB_i)}{P(A)} = \frac{P(B_i) \times P(A/B_i)}{\sum_{i=1}^k P(B_i) \times P(A/B_i)}$$

This last equation is often called Law of Total Probability.

Example 1

A rare genetic disease (occurring in 1 out of 1000 people) is diagnosed using a DNA screening test. The test has false positive rate of 0.5%, meaning that $P(\text{test positive} / \text{no disease}) = 0.005$. Given that a person has tested positive, what is the probability that this person actually has the disease? First, guess the answer, then read on.

Solution

Let's reason in terms of actual numbers of people, for a change.

Imagine 1000 people, 1 of them having the disease. How many out of 1000 will test positive? One that actually has the disease, and about 5 disease-free people who would test false positive. Thus, $P(\text{disease/test positive}) \approx \frac{1}{6}$.

It is left as an exercise for the reader to write down the formal probability calculation.

Example 2

At a certain assembly plant, three machines make 30%, 45%, and 25%, respectively, of the products. It is known from the past experience that 2%, 3% and 2% of the products made by each machine, respectively, are defective. Now, suppose that a finished product is randomly selected.

- What is the probability that it is defective?
- If a product were chosen randomly and found to be defective, what is the probability that it was made by machine 3?

Solution

Consider the following events:

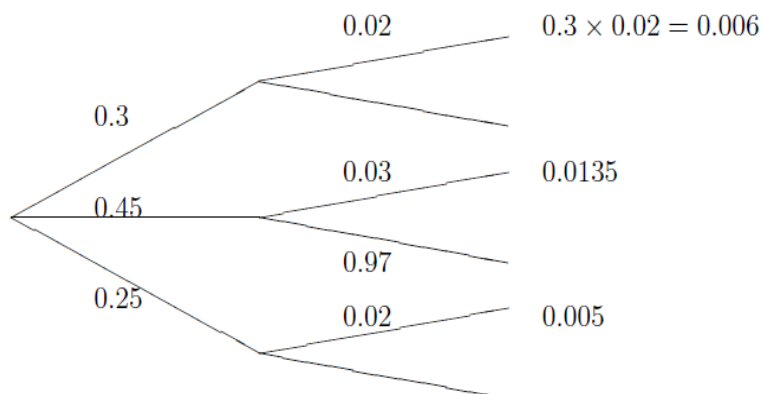
A : the product is defective and B_i : the product is made by machine $i=1, 2, 3$,

Applying additive and multiplicative rules, we can write

$$(a) \quad P(A) = P(B_1) \times P(A/B_1) + P(B_2) \times P(A/B_2) + P(B_3) \times P(A/B_3) \\ = (0.3)(0.02) + (0.45)(0.03) + (0.25)(0.02) = 0.006 + 0.0135 + 0.005 = 0.0245$$

$$(b) \quad \text{Using Bayes' rule } P(B_3/A) = \frac{P(B_3) \times P(A/B_3)}{P(A)} = \frac{0.005}{0.0245} = 0.2041$$

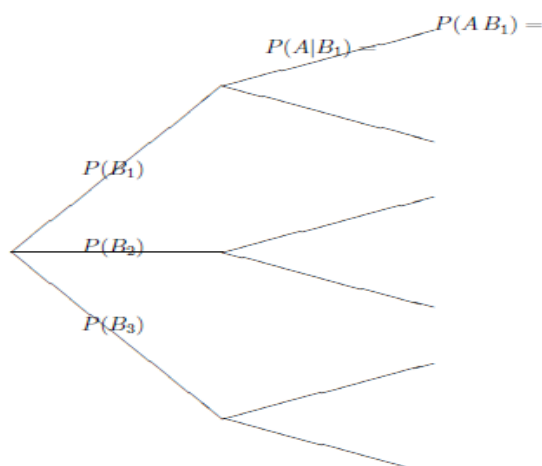
This calculation can also be represented using a tree diagram as follows



Here, the first branching represents probabilities of the events B_i and the second branching represents conditional probabilities $P(A/B_i)$. The probabilities of intersections, given by the products, are on the right. $P(A)$ is their sum.

Exercises

- 1) Lucy is undecided as to whether to take a Math course or a Chemistry course. She estimates that her probability of receiving an A grade would be 0.5 in a math course, and $\frac{2}{3}$ in a chemistry course. If Lucy decides to base her decision on the flip of a fair coin, what is the probability that she gets an A?
- 2) Of the customers at a gas station, 70% use regular gas, and 30% use diesel. Of the customers who use regular gas, 60% will fill the tank completely, and of those who use diesel, 80% will fill the tank completely.
 - a) What percent of all customers will fill the tank completely?
 - b) If a customer has filled up completely, what is the probability it was a customer buying diesel?
- 3) In 2004, 57% of White households directly and/or indirectly owned stocks, compared to 26% of Black households and 19% of Hispanic households. The data for Asian households is not given, but let's assume the same rate as for Whites. Additionally, 77% of households are classified as either White or Asian, 12% as African American, and 11% as Hispanic.
 - a) What proportion of all families owned stocks?
 - b) If a family owned stock, what is the probability it was White/Asian?
- 4) Drawer one has five pairs of white and three pairs of red socks, while drawer two has three pairs of white and seven pairs of red socks. One drawer is selected at random a pair of socks is selected at random from that drawer.
 - a) What is the probability that it is a white pair of socks?
 - b) Suppose a white pair of socks is obtained. What is the probability that it came from drawer two?
- 5) For an on-line electronics retailer, 5% of customers who buy Zony digital cameras will return them, 3% of customers who buy Lucky Star digital cameras will return them, and 8% of customers who buy any other brand will return them. Also, among all digital cameras bought, there are 20% Zony's and 30% Lucky Stars. Fill in the tree diagram and answer the questions.
 - a) What percent of all cameras are returned?
 - b) If the camera was just returned, what is the probability it is a Lucky Star?
 - c) What percent of all cameras sold were Zony and were not returned?



- 6) Three newspapers, A, B, and C are published in a certain city. It is estimated from a survey that that of the adult population: 20% read A, 16% read B, 14% read C, 8% read both A and B, 5% read both A and C, 4% read both B and C, 2% read all three. What percentage reads at least one of the papers? Of those that read at least one, what percentage reads both A and B?
- 7) Suppose $P(A/B) = 0.3$, $P(B) = 0.4$ and $P(B/A) = 0.6$. Find $P(A)$ and $P(A \cup B)$
- 8) This is the famous Monty Hall problem. A contestant on a game show is asked to choose among 3 doors. There is a prize behind one door and nothing behind the other two. You (the contestant) have chosen one door. Then, the host is flinging one other door open, and there's nothing behind it. What is the best strategy? Should you switch to the remaining door, or just stay with the door you have chosen? What is your probability of success (getting the prize) for either strategy?
- 9) There are two children in a family. We overheard about one of them referred to as a boy.
 - a) Find the probability that there are 2 boys in the family.
 - b) Suppose that the oldest child is a boy. Again, find the probability that there are 2 boys in the family. [Why is it different from part (a)?]
- 10) At a university, two students were doing well for the entire semester but failed to show up for a final exam. Their excuse was that they travelled out of state and had a flat tire. The professor gave them the exam in separate rooms, with one question worth 95 points: "which tire was it?". Find the probability that both students mentioned the same tire.
- 11) In firing the company's CEO, the argument was that during the six years of her tenure, for the last three years the company's market share was lower than for the first three years. The CEO claims bad luck. Find the probability that, given six random numbers, the last three are the lowest among six.

6 DISCRETE PROBABILITY DISTRIBUTIONS

In this section, we will consider random quantities that are usually called random variables.

Introduction

In application of probability, we are often interested in a number associated with the outcome of a random experiment. Such a quantity whose value is determined by the outcome of a random experiment is called a **random variable**. It can also be defined as any quantity or attribute whose value varies from one unit of the population to another.

A **discrete** random variable is a function whose range is finite and/or countable, i.e. it can only assume values in a finite or countably infinite set of values. A **continuous** random variable is one that can take any value in an interval of real numbers. (There are *uncountably* many real numbers in an interval of positive length.)

6.1 Discrete Random Variables

A random variable X is said to be discrete if it can take on only a finite or countable number of possible values x . Consider the experiment of flipping a fair coin three times. The number of tails that appear is noted as a discrete random variable. $X =$ "number of tails that appear in 3 flips of a fair coin". There are 8 possible outcomes of the experiment: namely the sample space consists of

$$S = \{HHH, HHT, HTH, HTT, THH, THT, TTH, TTT\}$$

$$X = \{0, 1, 1, 2, 1, 2, 2, 3\}$$

are the corresponding values taken by the random variable X .

Now, what are the possible values that X takes on and what are the probabilities of X taking a particular value?

From the above we see that the possible values of X are the 4 values

$$X = \{0, 1, 2, 3\}$$

ie the sample space is a disjoint union of the 4 events $\{X = j\}$ for $j=0,1,2,3$

Specifically in our example :

$$\{X = 0\} = \{HHH\}$$

$$\{X = 1\} = \{HHT, HTH, THH\}$$

$$\{X = 2\} = \{TTH, HTT, THT\}$$

$$\{X = 3\} = \{TTT\}$$

Since for a fair coin we assume that each element of the sample space is equally likely (with probability $\frac{1}{8}$), we find that the probabilities for the various values of X , called the *probability distribution* of X or the *probability mass function (pmf)*, can be summarized in the following table listing the possible values beside the probability of that value

x	0	1	2	3
P(X=x)	$\frac{1}{8}$	$\frac{3}{8}$	$\frac{3}{8}$	$\frac{1}{8}$

Note: The probability that X takes on the value x , ie $p(X = x)$, is defined as the sum of the probabilities of all points in S that are assigned the value x .

We can say that this pmf places mass $\frac{3}{8}$ on the value $X = 2$.

The “masses” (or probabilities) for a pmf should be between 0 and 1.

The total mass (i.e. total probability) must add up to 1.

Definition: The **probability mass function** of a discrete variable is a, table, formula or graph that specifies the proportion (or probabilities) associated with each possible value the random variable can take. The mass function $P(X = x)$ (or just $p(x)$) has the following properties:

$$0 \leq p(x) \leq 1 \text{ and } \sum_{\text{all } x} p(x) = 1$$

More generally, let X have the following properties

- It is a discrete variable that can only assume values x_1, x_2, \dots, x_n
- The probabilities associated with these values are $P(X = x_1) = p_1$, $P(X = x_2) = p_2, \dots, P(X = x_n) = p_n$

Then X is a discrete random variable if $0 \leq p_i \leq 1$ and $\sum_{i=1}^n p_i = 1$

Remark: We denote random variables with capital letters while realized or particular values are denoted by lower case letters.

Example 1

Two tetrahedral dice are rolled together once and the sum of the scores facing down was noted. Find the pmf of the random variable ‘the sum of the scores facing down.’

Solution

+	1	2	3	4
1	2	3	4	5
2	3	4	5	6
3	4	5	6	7
4	5	6	7	8

$$X = \{1, 2, 3, 4, 5, 6, 7, 8\}$$

Therefore the pmf is given by the table below

x	2	3	4	5	6	7	8
P(X=x)	$\frac{1}{16}$	$\frac{1}{8}$	$\frac{3}{16}$	$\frac{1}{4}$	$\frac{3}{16}$	$\frac{1}{8}$	$\frac{1}{16}$

This can also be written as a function

$$P(X = x) = \begin{cases} \frac{x-1}{16} & \text{for } x = 2, 3, 4, 5 \\ \frac{9-x}{16} & \text{for } x = 6, 7, 8 \end{cases}$$

Example 2

The pmf of a discrete random variable W is given by the table below

w	-3	-2	-1	0	1
P(W=w)	0.1	0.25	0.3	0.15	d

Find the value of the constant d, $P(-3 \leq w < 0)$, $P(w > -1)$ and $P(-1 < w < 1)$

Solution

$$\sum_{\text{all } w} p(W = w) = 1 \Rightarrow 0.1 + 0.25 + 0.3 + 0.15 + d = 1 \Rightarrow d = 0.2$$

$$P(-3 \leq w < 0) = P(W = -3) + P(W = -2) + P(W = -1) = 0.65$$

$$P(w > -1) = P(w = 0) + P(w = 1) = 0.15 + 0.2 = 0.35$$

$$P(-1 < w < 1) = P(W = 0) = 0.15$$

Example 3

A discrete random variable Y has a pmf given by the table below

y	0	1	2	3	4
P(Y=y)	c	2c	5c	10c	17c

Find the value of the constant c hence computes $P(1 \leq Y < 3)$

Solution

$$\sum_{\text{ally}} p(Y = y) = 1 \Rightarrow c(1 + 2 + 5 + 10 + 17) = 1 \Rightarrow c = \frac{1}{35}$$

$$P(1 \leq Y < 3) = P(Y = 1) + P(Y = 2) = \frac{2}{35} + \frac{5}{35} = \frac{1}{5}$$

Exercise

1. A die is loaded such that the probability of a face showing up is proportional to the face number. Determine the probability of each sample point.
2. Roll a fair die and let X be the square of the score that show up. Write down the probability distribution of X hence compute $P(X < 15)$ and $P(3 \leq X < 30)$
3. Let X be the random variable the number of fours observed when two dice are rolled together once. Show that X is a discrete random variable.
4. The pmf of a discrete random variable X is given by $P(X = x) = kx$ for $x = 1, 2, 3, 4, 5, 6$
Find the value of the constant k, $P(X < 4)$ and $P(3 \leq X < 6)$
5. A fair coin is flip until a head appears. Let N represent the number of tosses required to realize a head. Find the pmf of N c , $P(N < 2)$ and $P(N \geq 2)$
6. A discrete random variable Y has a pmf given by $P(Y = y) = c\left(\frac{3}{4}\right)^y$ for $y = 0, 1, 2, \dots$
Find the value of the constant c , $P(X < 3)$ and $P(X \geq 3)$
7. Verify that $f(x) = \frac{2x}{k(k+1)}$ for $y = 0, 1, 2, \dots, k$ can serve as a pmf of a random variable X.
8. For each of the following determine c so that the function can serve as a pmf of a random variable X.
 - a. $f(x) = c$ for $x = 1, 2, 3, 4, 5$
 - b. $f(x) = cx$ for $x = 1, 2, 3, 4, 5$
 - c. $f(x) = cx^2$ for $x = 0, 1, 2, \dots, k$
 - d. $f(x) = \frac{c}{2}$ for $x = -1, 0, 1, 2$
 - e. $f(x) = \frac{(x-2)}{c}$ for $x = 1, 2, 3, 4, 5$
 - f. $f(x) = \frac{(x^2 - x + 1)}{c}$ for $x = 1, 2, 3, 4, 5$
 - g. $f(x) = c(x^2 + 1)$ for $x = 0, 1, 2, 3$
 - h. g) $f(x) = cx({}_3C_x)$ for $x = 1, 2, 3$
 - i. $f(x) = c\left(\frac{1}{6}\right)^x$ for $x = 0, 1, 2, 3, \dots$
 - j. $f(x) = c2^{-x}$ for $x = \text{for } x = 0, 1, 2, \dots$
9. A coin is loaded so that heads is three times as likely as the tails.

- a. For 3 independent tosses of the coin find the pmf of the total number of heads realized and the probability of realizing at most 2 heads.
 - b. A game is played such that you earn 2 points for a head and loss 5 points for a tail. Write down the probability distribution of the total scores after 4 independent tosses of the coin
10. For an on-line electronics retailer, X = “the number of Zony digital cameras returned per day” follows the distribution given by
- | | | | | | | |
|----------|------|-----|-----|-----|------|-----|
| x | 0 | 1 | 2 | 3 | 4 | 5 |
| $P(X=x)$ | 0.05 | 0.1 | t | 0.2 | 0.25 | 0.1 |
- Find the value of t and $P(X > 3)$
11. Out of 5 components, 3 are domestic and 2 are imported. 3 components are selected at random (without replacement). Obtain the PMF of X = “number of domestic components picked” (make a table).

6.2 Expectation and Variance of a Random Variable

6.2.2 Expected Values

One of the most important things we'd like to know about a random variable is: what value does it take on average? What is the average price of a computer? What is the average value of a number that rolls on a die? The value is found as the average of all possible values, weighted by how often they occur (i.e. probability)

Definition: Let X be a discrete r.v. with probability function $p(x)$. Then the **expected value** of

X , denoted $E(X)$ or μ , is given by $E(x) = \mu = \sum_{x=-\infty}^{\infty} xp(X = x)$.

Theorem: Let X be a discrete r.v. with probability function $p(X=x)$ and let $g(x)$ be a real-valued function of X . ie $g: \mathbb{R} \rightarrow \mathbb{R}$, then the expected value of $g(x)$ is given by

$$E[g(x)] = \sum_{x=-\infty}^{\infty} g(x)p(X = x).$$

Proof (left as exercise)

Theorem: Let X be a discrete r.v. with probability function $p(x)$. Then

- (i) $E(c) = c$, where c is any real constant;
- (ii) $E[ax + b] = a\mu + b$ where a and b are constants
- (iii) $E[kg(x)] = kE[g(x)]$ where $g(x)$ is a real-valued function of X
- (iv) $E[ag_1(x) \pm bg_2(x)] = aE[g_1(x)] \pm bE[g_2(x)]$ and in general $E\left[\sum_{i=1}^n c_i g_i(x)\right] = \sum_{i=1}^n c_i E[g_i(x)]$

where $g_i(x)$ are real-valued functions of X .

This property of expectation is called *linearity property*

Proof

- (i) $E[c] = \sum_{all\ x} cP(X = x) = c \sum_{all\ x} P(X = x) = c(1) = c$
- (ii) $E[ax + b] = \sum_{all\ x} (ax + b)P(x) = \sum_{all\ x} axP(x) + \sum_{all\ x} bP(x) = a \sum_{all\ x} xP(x) + b \sum_{all\ x} P(x) = a\mu + b$
- (iii) $E[kg(x)] = \sum_{all\ x} kg(x)P(X = x) = k \sum_{all\ x} g(x)P(X = x) = kE[g(x)]$

(iv) $E[ag_1(x) \pm bg_2(x)] = E[ag_1(x)] \pm E[bg_2(x)] = aE[g_1(x)] \pm bE[g_2(x)]$ from part iii

6.2.3 Variance and Standard Deviation

Definition: Let X be a r.v with mean $E(X) = \mu$, the **variance** of X , denoted σ^2 or $\text{Var}(X)$, is given by $\text{Var}(X) = \sigma^2 = E(X - \mu)^2$. The units for variance are square units. The quantity that has the correct units is **standard deviation**, denoted σ . It's actually the positive square root of $\text{Var}(X)$.

$$\sigma = \sqrt{\text{Var}(X)} = \sqrt{E(X - \mu)^2}.$$

Theorem: $\text{Var}(X) = E(X - \mu)^2 = E(X)^2 - \mu^2$

Proof:

$$\text{Var}(X) = E(X - \mu)^2 = E(X^2 - 2X\mu + \mu^2) = E(X)^2 - 2\mu E(X) + \mu^2 = E(X)^2 - \mu^2 \text{ Since } E(X) = \mu$$

Theorem: $\text{Var}(aX + b) = a^2 \text{var}(X)$

Proof:

Recall that $E[aX + b] = a\mu + b$ therefore

$$\text{Var}(aX + b) = E[(aX + b) - (a\mu + b)]^2 = E[a(X - \mu)]^2 = E[a^2(X - \mu)^2] = a^2 E[(X - \mu)^2] = a^2 \text{var}(X)$$

Remark

(i) The expected value of X always lies between the smallest and largest values of X .

(ii) In computations, bear in mind that variance cannot be negative!

Example 1

Given a probability distribution of X as below, find the mean and standard deviation of X .

x	0	1	2	3
P(X=x)	1/8	1/4	3/8	1/4

Solution

x	0	1	2	3	total
$p(X = x)$	1/8	1/4	3/8	1/4	1
$xp(X = x)$	0	1/4	3/4	3/4	7/4
$x^2 p(X = x)$	0	1/4	3/2	9/4	4

$$E(X) = \mu = \sum_{x=0}^3 xp(X = x) = 1.75 \text{ and}$$

standard deviation

$$\sigma = \sqrt{E(X^2) - \mu^2} = \sqrt{4 - 1.75^2} \approx 0.968246$$

Example 2

The probability distribution of a r.v X is as shown below, find the mean and standard deviation of; a) X b) $Y = 12X + 6$.

x	0	1	2
P(X=x)	1/6	1/2	1/3

Solution

x	0	1	2	total
$p(X = x)$	1/6	1/2	1/3	1
$xp(X = x)$	0	1/2	2/3	7/6
$x^2 p(X = x)$	0	1/2	4/3	11/6

$$E(X) = \mu = \sum_{x=0}^2 xp(X = x) = 7/6 \text{ and}$$

$$E(X^2) = \sum_{x=0}^2 x^2 p(X = x) = 11/6$$

$$\text{Standard deviation } \sigma = \sqrt{E(X^2) - \mu^2} = \sqrt{11/6 - (7/6)^2} = \sqrt{17/6} \approx 1.6833$$

Now $E(Y) = 12E(X) + 6 = 12(\frac{7}{6}) + 6 = 20$

$Var(Y) = Var(12X + 6) = 12^2 \times Var(X) = 144 \times \sqrt{\frac{17}{6}} \approx 242.38812$

Exercise

1. Suppose X has a probability mass function given by the table below

x	2	3	4	5	6
P(X=x)	0.01	0.25	0.4	0.3	0.04

Find the mean and variance of; X

2. Suppose X has a probability mass function given by the table below

x	11	12	13	14	15
P(X=x)	0.4	0.2	0.2	0.1	0.1

Find the mean and variance of; X

3. Let X be a random variable with $P(X = 1) = 0.2$, $P(X = 2) = 0.3$, and $P(X = 3) = 0.5$. What is the expected value and standard deviation of; a) X b) $Y = 5X - 10$?
4. A random variable W has the probability distribution shown below,

w	0	1	2	3
P(W=w)	2d	0.3	d	0.1

Find the values of the constant d hence determine the mean and variance of W. Also find the mean and variance of $Y = 10X + 25$

5. A random variable X has the probability distribution shown below,

x	1	2	3	4	5
P(X=x)	7c	5c	4c	3c	c

Find the values of the constant c hence determine the mean and variance of X.

6. The random variable Z has the probability distribution shown below,

z	2	3	5	7	11
P(Z=z)	$\frac{1}{6}$	$\frac{1}{3}$	$\frac{1}{4}$	x	y

If $E(Z) = 4\frac{2}{3}$, find the values of x and y hence determine the variance of Z

7. A discrete random variable M has the probability distribution $f(m) = \begin{cases} \frac{m}{36}, & m = 1, 2, 3, \dots, 8 \\ 0, & elsewhere \end{cases}$,
find the mean and variance of M

8. For a discrete random variable Y the probability distribution is $f(y) = \begin{cases} \frac{5-y}{10}, & y = 1, 2, 3, 4 \\ 0, & elsewhere \end{cases}$,
calculate $E(Y)$ and $var(Y)$

9. Suppose X has a pmf given by $f(x) = \begin{cases} kx & \text{for } x = 1, 2, 3, 4 \\ 0, & elsewhere \end{cases}$, find the value of the constant k

hence obtain the mean and variance of X

10. A team of 3 is to be chosen from 4 girl and 6 boys. If X is the number of girls in the team, find the probability distribution of X hence determine the mean and variance of X
11. A fair six sided die has; '1' on one face, '2' on two of its faces and '3' on the remaining three faces. The die is rolled twice. If T is the total score write down the probability distribution of T hence determine;
- the probability that T is more than 4
 - the mean and variance of T