



MIS171 - Summary Business Analytics

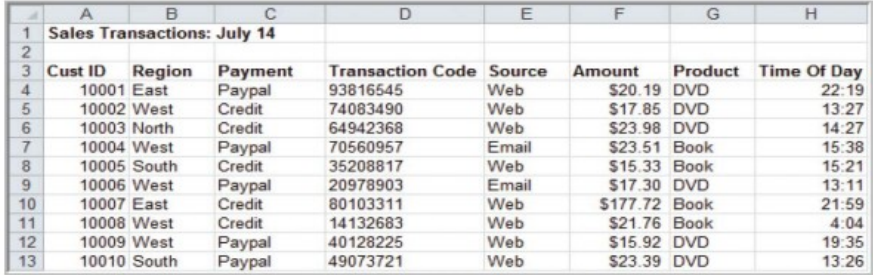
Business Analytics (Deakin University)

	Word	Definition	Notes
Week 1 (Topic 1) - Chapter ONE (Role of Business Analytics in real world context)			
Business Analytics	Business Analytics	<p>Process of transforming data into actions through analysis and insights in the context of organisational decision making and problem solving</p> <p>Using data, IT, statistical analysis, quantitative methods, and mathematical or computer-based models to help managers gain an improved insight about their business operations and make better fact-based decisions</p> <ul style="list-style-type: none"> • Report using historical info • Gives info to enable future predictions • End goal: add value through insight + turn data into info • Makes distinction b/w relevant + irrelevant knowledge 	
	Purpose of BA	<ul style="list-style-type: none"> • ID-ing valuable data re business' strategy + objectives • Internal value = revenue growth 	
	Importance of BA	<ul style="list-style-type: none"> • aids decision making • decisions made using analysis are better than those made through gut instinct • use of analytics = profitability and revenue 	
	Adds value by its:	<ul style="list-style-type: none"> • Business relevancy • Actionable insight • Performance + value measurement 	
	BA helps understand:	<ul style="list-style-type: none"> • What will happen • Why happened • Best course of action 	
	BA applications	<ul style="list-style-type: none"> • mgmt. of customer rels • financial + marketing activities • supply chain mgmt. • HR planning • Pricing decisions 	
Scope of BA	Descriptive analysis	<p>want to know about the past</p> <ul style="list-style-type: none"> • Most commonly used + most well understood type of analytics • Use data to understand past and present performance to make important decisions • Summarizes data into meaningful charts and reports <p>Focuses on:</p> <ul style="list-style-type: none"> • descriptive measures • data visualisation • probability distributions / sample + estimation • statistical inference 	
Scope of	Predictive	want to know abc	

BA	analysis	<ul style="list-style-type: none"> Analyses past performance in an effort to predict the future by examining historical data, detecting patterns or relationships in data <p>Techniques include:</p> <ul style="list-style-type: none"> regression forecasting 	
Scope of BA	Prescriptive analysis	<p>making decisions / optimisation</p> <ul style="list-style-type: none"> Uses optimization to identify the best alternative to minimize or maximise some objective Addresses questions such as: <ul style="list-style-type: none"> How much should we produce to maximize profit? What is the best way of shipping goods from our factory to minimize costs? <p>Techniques include:</p> <ul style="list-style-type: none"> optimisation simulation 	
	BA V analytics	<p>Business Analytics</p> <ul style="list-style-type: none"> BA = making info have contextual relevancy + delivering real value Insight that is actionable to create value <p>Analytics</p> <ul style="list-style-type: none"> A = finding interesting things in lrg amounts of data Focuses on creation of insight Simply answers question without necessary return <p>80% accuracy <u>with</u> actionable insight better than 98% accuracy <u>without</u> actionable insight</p>	
	Examples of analytics	<ul style="list-style-type: none"> Reporting: summary of historical data Trending: ID-ing patters Segmentation: ID-ing similarities in data Predictive modelling: using historical data to predict future events 	
	Mgmt.	<p>Process which org achieves goals through use of resources</p> <p>(ppl, money, materials + info = input, achieving goals = output)</p>	
	Role of manager	<ul style="list-style-type: none"> Interpersonal role: figurehead, leader, liaison Informational role: monitor, disseminator, analyser, spokesperson Decisional role: entrepreneur, disturbance handler, resource allocator, negotiator 	
	Decision	Choice b/w alternatives made by individual/group	
	Nature of decisions	<ul style="list-style-type: none"> operational control = executing specific tasks efficiently + effectively mgmt. control = acquiring + using resources efficiently strategic planning = long term goals + policies for growth + resource allocation 	
	Difficulties in	<ul style="list-style-type: none"> Number of alternatives 	

	decision making	<ul style="list-style-type: none"> • Time pressures • Conduct analysis • May be necessary to access info/consult with expert = therefore computerised analysis streamlines decision making	
	Process of decision process	1. Intelligence: What is the problem? 2. Design: what are my options? 3. Choice: pick an option + decide how to implement (if choice doesn't work/achieve goal = start process again)	
Type of Decision	Structured decisions	Std problems <ul style="list-style-type: none"> • routine + repetitive problems w std solutions eg. Order entry, accounts receivable	
Type of Decision	Unstructured decisions	Complex problems <ul style="list-style-type: none"> • no cut-dry solution eg. Planning new service offering, hiring executive	
Type of Decision	Semi-structured decisions	Combination of std solution procedures + individual judgement eg. Evaluating employees, setting marketing budgets	
Business Intelligence	Business Intelligence	decision support applications + technologies + processes	
	BI applications	provide users with view of what has happened	
	Data mining	process of searching for valuable info in lrg database, warehousing or mart <ul style="list-style-type: none"> • helps explain why it is happening + predicts what will happen in future • predicts trends + behaviours • identifies prev. unknown patterns • able to provide predictive info eg. For targeted marketing, forecasting bankruptcy/defaults eg. can use prev. promotional info to identify ppl most likely to response to similar offers	
	Drill down	ability to go into more detail	
	Models	simplified representations or abstracts of reality	
Decision support system	Decision support system (DSS)	<ul style="list-style-type: none"> • combines models + data to analyse semi-structured + unstructured problems • enables ppl access to data to manipulate data + conduct appropriate analyses • enhances learning + contribute to all levels of decision making 	
Type of DSS	Sensitivity analysis	study of the impact one or more parts of decision has on other parts <ul style="list-style-type: none"> • examines impact of input variable changes on output variables • enables system to adapt to changing conditions + varying requirements of diff situations • provide better understanding of model + problem model describes 	

Type of DSS	What-if analysis	attempts to predict impact of changes based on assumptions (input data) on proposed solution	
Type of DSS	Goal-seeking analysis	<ul style="list-style-type: none"> find the value of inputs necessary to achieve desired level of output finding out what is needed to achieve certain goal 	
	Dashboards	<ul style="list-style-type: none"> evolved from executive info systems = designed for info needs of execs provides access to timely info + direct access to mgmt. reports enable managers to examine reports + drill down into detailed info 	
	Data visualisation technologies	presenting data in the form of graphs or tables to enable users to more easily process and understand data	
	Geographic Info Systems (GIS)	system for capturing, integrating, manipulating + displaying data using digitised maps	
	Geo-coding	identifying geographical location of every digital record <ul style="list-style-type: none"> enables users to generate info for planning, problem solving + decision making graphical format makes it easy for managers to visualise data 	
	Reality mining	able to extract info from usage patterns of mobile phones + other wireless devices	
	Real-time BI	use of data as and when it happens <ul style="list-style-type: none"> better ways to communication and make decisions that affect customers 	
	Corp performance mgmt. (CPM)	<ul style="list-style-type: none"> monitoring + managing performance in accord with KPIs BI allows ppl to view info + insights re co. KPIs 	
Big Data	Big Data	massive amounts of data / data sets	
	Volume	amount of data produce	
	Variety	different forms EG. msgs, updates + images posted to social media, readings sensors, GPS signals etc	
	Velocity	how fast data is produced + how fast data needs to be processed to meet demand	
	Data	facts + figures collected	
	Information	from analysing data	
	Data matrix	dataset that is stored digitally in spreadsheet or similar	
Dataset	Dataset	data collected in particular study	
	Entities	people, places or things which we store + maintain info	
	variable / attribute	characteristic of interest of entity	

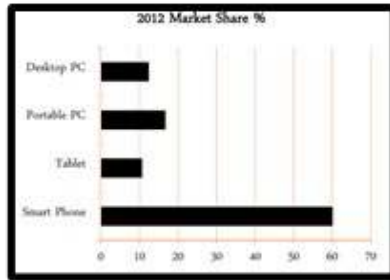

		 <p>Entities</p> <p>Variables or Attributes</p> <p>Records</p>	
	Record	set of measurements collected for particular entity is a record / observation	
	Data deluge	prevalence of: <ul style="list-style-type: none"> • automatic data collection • electronic instrumentation • online transaction processing • growing recognition of untapped value • recognition is driving development of BA 	

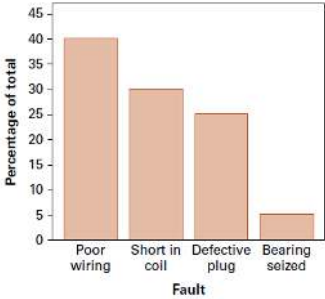
Week 2 (Topic 2) - Chapter TWO (What is Data?)			
Statistics	Statistics	study of variation in data relates to collection, analysis, interpretation & presentation of data statistical methods: <ul style="list-style-type: none"> • summarise collection of data • draw inferences about entire population • make predictions or forecasts 	
statistical studies	experimental studies	the variables of interest are first identified → one or more factors are controlled so data can be obtained re how it affects variables	
	observational studies	no attempt is made to control or influence variables	
	Descriptive statistics	tabular, graphical + numerical methods used to summarise data	
	inference statistics	process of using data to obtain from a sample to make estimates + test claims re characteristics of population uses sample data to reach conclusions re population from which sample drawn <div> Process of → Inference Statistics <pre> graph TD 1((1. Population consists of all employees at Conrobar. Average productivity is unknown)) --> 2[2. A random sample of 48 employees is examined] 2 --> 3[3. The sample data provide a sample average productivity of 98.9%] 3 --> 4[4. The sample average is used to estimate the population average] 4 --> 1 </pre> </div>	
Population	Population	Entire collection of objects (called units or subjects) of interest	
	Census	Collection of data on population	
	Sample	Subset of units in population <ul style="list-style-type: none"> • Representative of whole population • Sometimes this is only way to get info (eg. crash data on cars) 	
Types of data	Exploratory data analysis (EDA)	<ul style="list-style-type: none"> • first step • precursor to more formal + extensive analysis • numerical, tabular + graphical summaries are produced to summarise + highlight key aspects or special 	


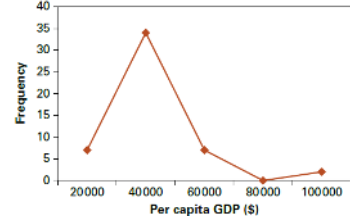
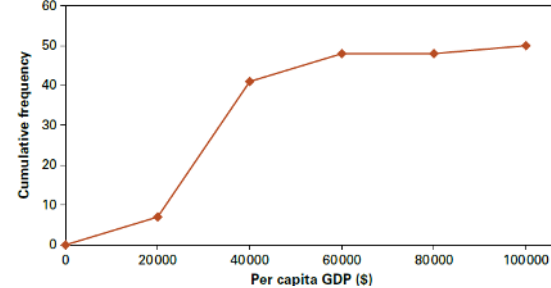
		features of data <ul style="list-style-type: none"> analysis sufficient for purpose of study 	
	Parameter	descriptive measure of population <ul style="list-style-type: none"> denoted by Greek letters <ul style="list-style-type: none"> μ = population mean σ^2 = population variance σ = population std deviation 	
	Statistic	descriptive measure of sample <ul style="list-style-type: none"> denoted by Roman letters <ul style="list-style-type: none"> \bar{x} = sample mean s^2 = sample variance s = sample standard deviation 	
	Parameter v Statistic	<ul style="list-style-type: none"> estimate value of parameter is desirable but impossible/not feasible b/c of time + \$\$ required for consensus instead use representative sample of pop + corresponding sample statistic to estimate parameter <p>EG. F&P wants to determine avg. no. of loads its 8kg washing machine can wash before repairs. Population = all 8kg washing machines Parameter = population mean = avg no. of washers per machine before repair Population mean = no. of washes for type of machine estimated from sample mean Statistician takes representative sample of machines, conducts trials + records no. of washes before repair for each machine then computes sample avg. no. of washes before repair.</p>	Research – don't quite understand
2 types of data	Qualitative data aka Categorical data	descriptive data subclass: nominal + ordinal <ul style="list-style-type: none"> Labels or names used to identify attributes of each entity Can be recorded in either numeric or nonnumeric formats EG. 'Yes or no', 'male or female' answers Usually counted or expressed as a portion or a percentage 	
	Quantitative data aka Numerical data	data that is expressed as number subclass: discrete + continuous <ul style="list-style-type: none"> Take numbers as their observed responses Numerical data can be converted to categorical data. EG Salary can be converted into low/medium/high. Cannot convert categorical data back to numerical data 	
	Discrete	measuring how many (whole numbers + not able to be fraction)	
	Continuous	measuring how much (decimal/fraction)	

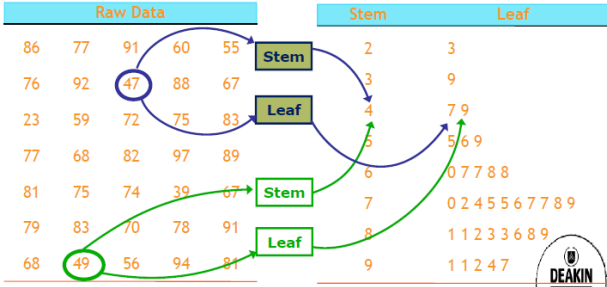
Scale of Measurement		<div> <div> SCALES OF MEASUREMENT - SUMMARY <p>Differences between measurements, true zero exists</p> <p>Differences between measurements but no true zero</p> <p>Ordered categories (rankings, order or scaling)</p> <p>Categories (no ordering or direction)</p> </div> <div> <p>Ratio Data</p> <p>Interval Data</p> <p>Ordinal Data</p> <p>Nominal Data</p> </div> <div> <p>Highest Level (Strongest forms of measurement)</p> <p>Lowest Level (Weakest forms of measurement)</p> </div> </div> <div> SCALES OF MEASUREMENT - SUMMARY <p>Numerical Data</p> <p>Categorical Data</p> <p>Ratio Data</p> <p>Interval Data</p> <p>Ordinal Data</p> <p>Nominal Data</p> <p>EXAMPLES:</p> <p>Height, age, monthly sales, delivery times</p> <p>Temperature in degrees Celsius, standardised exam score</p> <p>Service quality rating, student letter grades</p> <p>Marital status, customer's location, supplier's name</p> </div>	
	Nominal (categorical)	classification of categorical data that implies no ranking EG. male or female? green or blue?	
	Ordinal (categorical)	scale of measurement where values are assigned by ranking EG. rating customer service: worst - excellent	
	Interval (numerical)	ranking numerical data where differences are meaningful but there is no true zero point EG. shoe size, amount of liquid	
	Ratio (numerical)	ranking numerical data where differences involve a true zero point EG. length, weight, age salary	
2 broad types of data	Cross-sectional data	relates to group of items or individuals at given point in time	
	Time ordered (time series) data	relates to particular entity or situation at different points of time	
Data Sources	Primary	data collector is one using data EG. survey, experiment	
	Secondary	another org. or individual has collected data used for analysis by org or individual	
	Web data	enhanced with social media data	
	Data in Business	used in: <ul style="list-style-type: none"> • annual reports • accounting audits • profitability analysis • eco trends • market research 	

Week 3 (Topic 3) - Chapter THREE (Exploring Data - Data Visualisation)


Purpose of Data Visualisation	Analysis	purpose = insight - discovery, decision making + explanation <ul style="list-style-type: none">• able to comprehend huge amounts of data• important info from millions of data immediately visible and more comprehensible• provides insight into patterns + trends not readily visible• enables problems with data to become apparent																			
Categorical Data	Frequency table	gives frequency, proportion or % of value in each category / class <table border="1"><thead><tr><th>Type of device</th><th>2012 Shipments (in millions)</th><th>2012 Market Share</th></tr></thead><tbody><tr><td>Smart Phone</td><td>722.4</td><td>60.1%</td></tr><tr><td>Tablet</td><td>128.3</td><td>10.7%</td></tr><tr><td>Portable PC</td><td>202</td><td>16.8%</td></tr><tr><td>Desktop PC</td><td>148.4</td><td>12.4%</td></tr><tr><td>Total</td><td>1201.1</td><td>100%</td></tr></tbody></table>	Type of device	2012 Shipments (in millions)	2012 Market Share	Smart Phone	722.4	60.1%	Tablet	128.3	10.7%	Portable PC	202	16.8%	Desktop PC	148.4	12.4%	Total	1201.1	100%	
Type of device	2012 Shipments (in millions)	2012 Market Share																			
Smart Phone	722.4	60.1%																			
Tablet	128.3	10.7%																			
Portable PC	202	16.8%																			
Desktop PC	148.4	12.4%																			
Total	1201.1	100%																			
Categorical Data	Bar chart / Column chart	each category represented by a bar - the length indicates frequency, proportion or % of value in each category <div></div> <p>**use if comparing categories is most important</p>																			
Categorical Data	Pie graph	circle used to represent total - is divided into "slices", each representing category (used to observe proportion / market share) <div></div> <p>**use if observing portion of whole that lies in each category is most important</p>																			

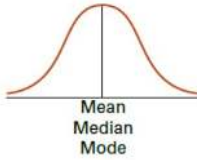
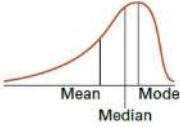
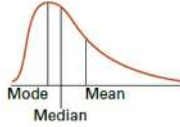
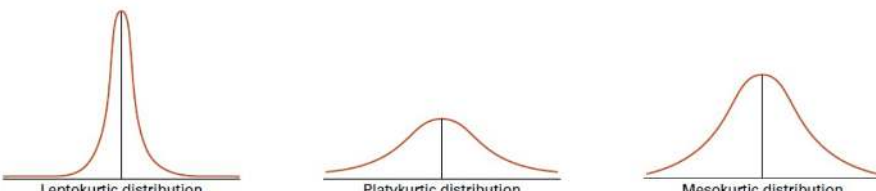
Categorical Data	Pareto chart	<p>vertical bar chart where no. + type are in order of magnitude (greatest to least)</p> <p>**used to <u>display causes of problems in products + processes</u></p>																																	
Numerical Data		<p>Common tabular + graphical techniques for organising + previewing include:</p> <ul style="list-style-type: none"> • arrays • frequency table / summary table • pie chart • bar chart • dot plots <p>useful tool to summarise:</p> <ul style="list-style-type: none"> • small dataset (<20 observations) • numerical data that is discrete + repeats frequently 																																	
Numerical Data	Arrays	putting data into ascending order																																	
Numerical Data	Frequency distribution	<p>summary table of data arranged into classes</p> <ul style="list-style-type: none"> • rules: <ul style="list-style-type: none"> - use 5 - 15 classes - class width = range/no. of classes - centre = mid point (used for graphing) • no. of observations in each class = frequency of class • allows quick visual interpretation of data • allows first look at shape of data <p>**use for large datasets + non-repeating values to summarise data</p>	<table border="1"> <thead> <tr> <th>Weekly Sales</th><th>Count</th><th>Percentage</th><th>Cum. Percentage</th></tr> </thead> <tbody> <tr> <td>0 kg < 200 kg</td><td>3</td><td>5.8%</td><td>5.8%</td></tr> <tr> <td>200 kg to < 400 kg</td><td>10</td><td>19.2%</td><td>25%</td></tr> <tr> <td>400 kg < 600 kg</td><td>16</td><td>30.8%</td><td>55.8%</td></tr> <tr> <td>600 kg < 800 kg</td><td>16</td><td>30.8%</td><td>86.6%</td></tr> <tr> <td>800 kg < 1000 kg</td><td>6</td><td>11.5%</td><td>98.1%</td></tr> <tr> <td>1000 kg < 1200 kg</td><td>1</td><td>1.9%</td><td>100%</td></tr> <tr> <td>Total</td><td></td><td>100%</td><td></td></tr> </tbody> </table>	Weekly Sales	Count	Percentage	Cum. Percentage	0 kg < 200 kg	3	5.8%	5.8%	200 kg to < 400 kg	10	19.2%	25%	400 kg < 600 kg	16	30.8%	55.8%	600 kg < 800 kg	16	30.8%	86.6%	800 kg < 1000 kg	6	11.5%	98.1%	1000 kg < 1200 kg	1	1.9%	100%	Total		100%	
Weekly Sales	Count	Percentage	Cum. Percentage																																
0 kg < 200 kg	3	5.8%	5.8%																																
200 kg to < 400 kg	10	19.2%	25%																																
400 kg < 600 kg	16	30.8%	55.8%																																
600 kg < 800 kg	16	30.8%	86.6%																																
800 kg < 1000 kg	6	11.5%	98.1%																																
1000 kg < 1200 kg	1	1.9%	100%																																
Total		100%																																	
	relative freq distribution	Equation: $\frac{\text{frequency in each class}}{\text{total no. of value}}$																																	
	% distribution	Equation: $\text{each relative frequency} \times 100$																																	
	cumulative % distribution	% of values that are less than certain value																																	

Numerical Data	Histogram	<p>Graphical representation of frequency, relative frequency, % distribution tables</p> <ul style="list-style-type: none"> Allows for a representation of the shape of the data set (skewness)  <p>Chocolates</p> <p>Frequency</p> <p>Weekly Sales (kg)</p> <p>**use for large datasets + non-repeating values to summarise data</p>	
Numerical Data	Frequency polygon	<ul style="list-style-type: none"> graph constructed by plotting dot for frequency at class mid points + connecting the dots  <p>Frequency</p> <p>Per capita GDP (\$)</p>	
Numerical Data	Ogive	<ul style="list-style-type: none"> cumulative frequency polygon plotted by graphing dot at each class endpoint for cumulative frequency value + connecting dots  <p>Cumulative frequency</p> <p>Per capita GDP (\$)</p>	
Numerical Data	Stem + leaf plots	<ul style="list-style-type: none"> another way to display continuous data major advantage is original data preserved displays info similar to histogram plot numbers constructed by separating each no. into 2 groups = stem + leaf leftmost digits are stems rightmost are leaves 	

			
Purpose of Data Visualisation	Communication	<ul style="list-style-type: none"> • tell story or show pattern that has already been discovered in data • focuses on the msg - clear + easy to understand • focus is on design of visualisation = should be honest, unambiguous + effective presentation of data 	
	Tufte's Principles of graphical display	<p>Graphical displays should:</p> <ul style="list-style-type: none"> • show the data • tell the truth (avoid distorting what the data has to say) • focus on the content (help the viewer think about the information rather than design) • encourage the eye to compare the data • make large data sets coherent • reveal data at several levels of detail • Closely integrate statistical and verbal descriptions 	

Week 4 (Topic 4) - Chapter FOUR (Numerical Summaries - Exploring Relationships)			
Summary measures	summary measures	<p>Single figure which attempts to summarise particular feature of set of data</p> <ul style="list-style-type: none"> if measures are computed for data from population = parameters if measures are computed for data from sample = statistics sample statistics is a point estimator of population parameter 	
What looking for in summary measures?	Measure of Central tendency	<p>Are averages:</p> <ul style="list-style-type: none"> mean: common average (arithmetic mean) median: middle value mode: most common value 	
	Median	<p>middle value of data set arranged in ascending order</p> <ul style="list-style-type: none"> is a measure of location – often reported for annual income + property value data if data set has any extreme values, median is preferred measure of central tendency extreme values can inflate the mean making it less typical 	
	Mode	<p>value that occurs with greatest frequency</p> <ul style="list-style-type: none"> greatest frequency can occur at 2+ different values if data has exactly 2 modes = data is bimodal if data has 2+ modes = multimodal if data has 1 mode = nominal 	
What looking for in summary measures?	Measure of Location	<ul style="list-style-type: none"> Averages are measures of location – indicating middle or centre of data set More specific measures of location can be found from ascending array: min, max, quartiles, percentiles 	
	Percentiles	<ul style="list-style-type: none"> Tell us location of certain percentages of data EG. Top 10%, smallest 1% etc <p>Process:</p> <ol style="list-style-type: none"> Organise data in ascending array Calculate percentile location: $i = \frac{P}{100}(n)$ <ul style="list-style-type: none"> If i is a whole number, percentile is average of values at the i and (next to i) positions If i is the whole number, the percentile value is found by rounding i up to the whole no. + reporting value at this position 	

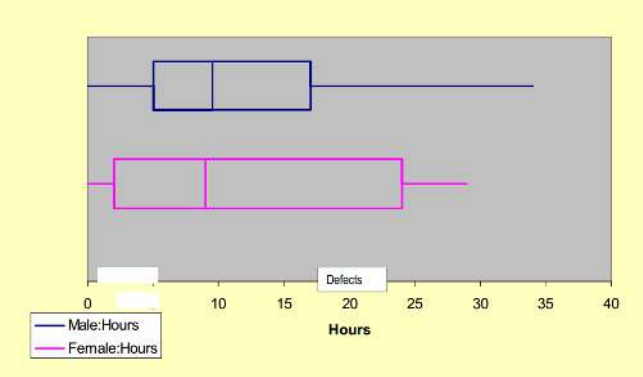
	Quartiles	<ul style="list-style-type: none"> Specific percentiles that are commonly used: <ul style="list-style-type: none"> 1st quartile: 25th percentile 2nd quartile: 50th percentile = median 3rd quartile: 75th percentile 	<p>Ascending array (customer waiting times in minutes)</p> <p>7 8 8 9 9 9 9 10 10 10 11 12 14 14 15</p> <p>Q1 or P25 $i = \frac{25}{100}(16)=4$ $Q_1 = \frac{9+9}{2} = 9$ minutes</p> <p>Q2 or P50 $i = \frac{50}{100}(16)=8$ $Q_2 = \frac{9+10}{2} = 9.5$ minutes</p> <p>Q3 or P75 $i = \frac{75}{100}(16)=12$ $Q_3 = \frac{11+12}{2} = 11.5$ minutes</p> 	
	Minimum	Lowest value		
	Maximum	Highest value		
What looking for in summary measures?	Measure of Variability	<p>Describes the spread or dispersion of data set</p> <ul style="list-style-type: none"> Distance measures: range, interquartile range Average variability: Standard deviation, and variance Relative variability: coefficient of variation 		
Distance measure	Range	<ul style="list-style-type: none"> Range of data set is difference b/w largest & smallest values Ignores all data points except 2 extreme ends of data set Very sensitive to smallest & largest data values 		
Distance measure	Inter Quartile Range (IQR)	<ul style="list-style-type: none"> Difference b/w 3rd quartile & 1st quartile $Q_1 - Q_3$ Range of middle 50% of data NOT sensitive to extreme data values 		
Average variation	Variation	Variation is expressed in squared units		
	Coefficient of Variation	<p>Indicates how large the std dev is in relation to the mean</p> <ul style="list-style-type: none"> Calculated as $\frac{std\ dev}{mean}$ Expressed at % Relative measure of variation <p>**Useful for comparing variability b/w data sets in different magnitudes or diff units</p>		
	Std Dev	<p>estimate of average deviation of value away from mean</p> <ul style="list-style-type: none"> Maintains original unit → preferred Popular measure of risk (esp. financial analysis) 		

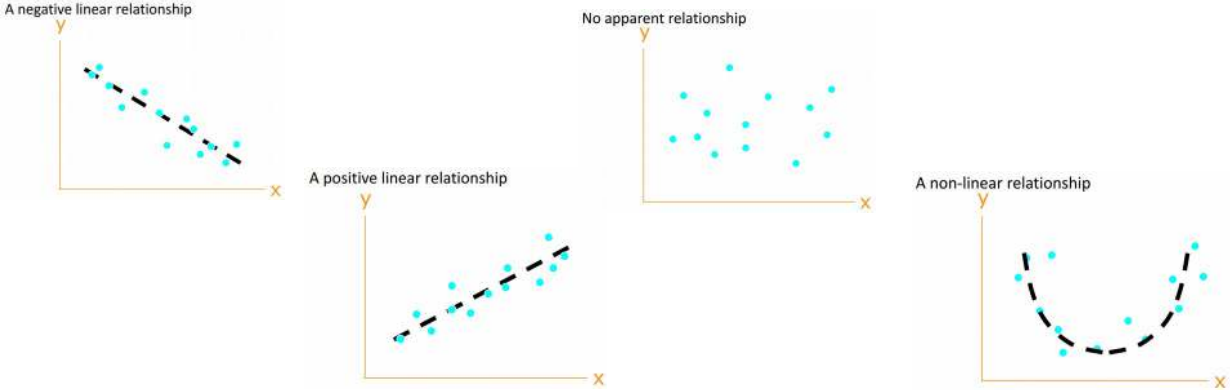
What looking for in summary measures?	Measure of Shape	<ul style="list-style-type: none"> Distribution of data where shape is symmetrical (bell curve)  <p>(a) Symmetrical distribution (no skewness)</p> <ul style="list-style-type: none"> Distribution of data where shape is skewed / asymmetric (lacks symmetry) <ul style="list-style-type: none"> Higher negative value = negatively skewed Higher positive value = positively skewed  <p>(b) Negatively skewed</p>  <p>(c) Positively skewed</p>	
	Kurtosis	Peakness of distribution:  <p>Leptokurtic distribution Platykurtic distribution Mesokurtic distribution</p>	
Relative location	Z scores	Relative measure of distance is an observation from the mean (re std dev) <ul style="list-style-type: none"> If z score is +3 or -3 = outlier <p>For example:</p> <ul style="list-style-type: none"> A dataset is normally distributed with a mean of 60 and standard deviation of 5. Determine the z score for a value of 70 $z = \frac{x - \mu}{\sigma}$ $z = \frac{70 - 60}{5} = 2.0$ <ul style="list-style-type: none"> a Z score of 2 means that 70 is 2 standard deviations above the mean 	
Relative location	Chebyshev's theorem	<ul style="list-style-type: none"> Applies to all distributions: <ul style="list-style-type: none"> At least 75% of the data values must be within Z=2 Standard deviations of the mean At least 89% of the data values must be within Z=3 Standard deviations of the mean At least 94% of the data values must be within Z=4 Standard deviations of the mean 	

		<ul style="list-style-type: none">At least $1 - 1/k^2$ values fall within $+ \text{ or } -k$Std dev of mean, regardless of shape of distribution <table><thead><tr><th>Number of Standard Deviations</th><th>Distance from the Mean</th><th>Minimum Proportion of Values Falling Within Distance</th></tr></thead><tbody><tr><td>$K = 2$</td><td>$\mu \pm 2 \sigma$</td><td>$1 - 1/2^2 = 0.75$</td></tr><tr><td>$K = 3$</td><td>$\mu \pm 3 \sigma$</td><td>$1 - 1/3^2 = 0.89$</td></tr><tr><td>$K = 4$</td><td>$\mu \pm 4 \sigma$</td><td>$1 - 1/4^2 = 0.94$</td></tr></tbody></table>	Number of Standard Deviations	Distance from the Mean	Minimum Proportion of Values Falling Within Distance	$K = 2$	$\mu \pm 2 \sigma$	$1 - 1/2^2 = 0.75$	$K = 3$	$\mu \pm 3 \sigma$	$1 - 1/3^2 = 0.89$	$K = 4$	$\mu \pm 4 \sigma$	$1 - 1/4^2 = 0.94$	
Number of Standard Deviations	Distance from the Mean	Minimum Proportion of Values Falling Within Distance													
$K = 2$	$\mu \pm 2 \sigma$	$1 - 1/2^2 = 0.75$													
$K = 3$	$\mu \pm 3 \sigma$	$1 - 1/3^2 = 0.89$													
$K = 4$	$\mu \pm 4 \sigma$	$1 - 1/4^2 = 0.94$													
Relative location	Empirical rule	<ul style="list-style-type: none">Guideline that states the approx. % of values that fall within given no. of std devs of mean = normally distributed (bell shaped)For data that is symmetrical / bell shaped:<ul style="list-style-type: none">Approx 68% of the data values lie within Z=1 Standard deviations of the meanApprox 95% of the data values lie within Z=2 Standard deviations of the meanApprox 99.7% of the data values lie within Z=3 Standard deviations of the mean													





What looking for in summary measures?	Outliers	<p>Data values that fall so far from average they are considered unusual</p> <ul style="list-style-type: none"> • Extreme extreme values • Tend to be separated from rest of data • Detecting outliers is important = need to be investigated more • May indicate: <ul style="list-style-type: none"> - Incorrectly recorded data value - Data value that was incorrectly included in data set - Potential problem - Potential opportunity • Detecting Outliers Graphically <ul style="list-style-type: none"> - Displaying data graphically = good way to spot potential outliers - Calculation rules are <u>not</u> the “be all and end all” of outlier detection - Various graphs can be used. Dot plots and box plots are often useful 	
Outliers	Empirical rule	<ul style="list-style-type: none"> • For symmetrical bell-shaped data • Further than 3 std dev from mean = potential outlier • Any data with z-score 3+ or -3 = potential outlier 	
Outliers	Tukey's 1.5 Step rule	<ul style="list-style-type: none"> • For non-bell shaped distributions • Works by calculating limits (fences) by determining quartiles + IQR <ul style="list-style-type: none"> - Lower fence = 1.5(IQR) below Q1 - Upper fence = 1.5(IQR) above Q3 • Value outside fences = potential outlier 	
	EXAMPLE: Tukey's rule	<ul style="list-style-type: none"> • Variable: Number of Daily iPhone Sales for 12 months for a Telco Company. • Summary Measures: <ul style="list-style-type: none"> - Min 13; - Q1 24; - Q3 46; - Max 97 (Positively Skewed Distribution) • Potential Outliers? <ul style="list-style-type: none"> - Lower fence: $Q1 - 1.5(IQR) = 24 - 1.5(22) = 24 - 33 = -9$ - Upper fence: $Q3 + 1.5(IQR) = 46 + 1.5(22) = 46 + 33 = 79$ • \therefore at least one potential outlier as 97 sales <u>does not</u> fall within the limits. 	

	Box + Whisker Plot	<ul style="list-style-type: none"> Five specific values are used: <ol style="list-style-type: none"> 1. Min value 2. 1st quartile 3. 2nd quartile / median 4. 3rd quartile 5. Max value Inner fences <ul style="list-style-type: none"> - IQR $Q_1 - Q_3$ - Lower inner fence = $Q_1 - 1.5\text{IQR}$ - Upper inner fence = $Q_3 + 1.5\text{IQR}$ Outer fence <ul style="list-style-type: none"> - Lower outer fence = $Q_1 - 3\text{IQR}$ - Upper outer fence = $Q_3 + 3.0\text{IQR}$ 		
What summary measure to use	Symmetrical distribution	Use: <ul style="list-style-type: none"> • Mean • Std dev 		
	Skewed distribution	Use: <ul style="list-style-type: none"> • Median (mean = too distorted) • IQR 		
Process: Relationship b/w 2 variables	Relationship b/w TWO variables	<ol style="list-style-type: none"> 1. Classify data (numerical / categorical) 2. Decide which variable is dependent variable & independent / explanatory variable 3. Select technique 4. Look for differences 		
	Techniques to use	<ul style="list-style-type: none"> • Numerical Dependent & Categorical Independent <ul style="list-style-type: none"> - Comparative summary measures - Multiple box plots • Categorical Dependent & Categorical Dependent <ul style="list-style-type: none"> - Cross tabulations (contingency tables) • Numerical Dependent & Numerical Independent <ul style="list-style-type: none"> - Scatter diagrams • Categorical Dependent & Numerical Independent <ul style="list-style-type: none"> - Numerical data can be converted to categorical data – cross tabulation 		
	Comparative summary measures	<ul style="list-style-type: none"> • Substantial differences = relationship • No or small differences = no relationship 		
	Multiple Box	<ul style="list-style-type: none"> • Useful way of displaying single numerical variable using 5 no. summary 		

	Plots	<ul style="list-style-type: none"> Places more than 1 box plot in single graph Enables us to quickly compare features + patterns across subgroups 	
	Cross tabulations	<ul style="list-style-type: none"> Displays results of 2 categorical variables together in 1 table Depend on type of cross tabulation, cells may contain: <ul style="list-style-type: none"> Absolute frequencies Relative frequencies Can be misleading if sample size of categories not similar Must calculate % % chosen to analyse should be where independent variable is located: <ul style="list-style-type: none"> If located in row = row % to compare If column = column % to compare <p>If % across row / column similar = no relationship b/w 2 categorical variables If % across rows dissimilar to columns = relationship b/w 2 categorical variables</p>	
	Relative frequencies constructed in 3 ways:	<ul style="list-style-type: none"> % of overall total % of row total % of column total 	
	Scatter diagram	<ul style="list-style-type: none"> Graphical presentation of reln b/w 2 numerical variables Independent variable shown on horizontal axis (x) Dependent variable shown on vertical axis (y) Pattern of plotted points suggest overall reln b/w variables Trend line = approximation of relationship 	

	Linear relationship	<ul style="list-style-type: none"> • closer points are to trend line = stronger relationship • measure direction + strength of relationship = correlation coefficient (r) • measure absolute strength of relationship = coefficient of determination (r^2) 	
	Causality	<ul style="list-style-type: none"> • Independent (explanatory) variable does not have causal effect on dependent variable • Cannot be demonstrated by data analysis techniques alone – must be through experiments where environment is controlled 	
SUMMARY		<ul style="list-style-type: none"> • More spread out the data: the larger the range + IQR + SD • The more concentrated or similar the data: the smaller the range + IQR + SD • If the values are the same: the range + IQR + SD will be zero • No measure of variation can ever be negative 	

Week 5 (Topic 5) - Chapter FIVE (Discrete Probability Distributions)

	Probability & Probability distributions	Enable us to develop models that take into account uncertainty <ul style="list-style-type: none">Decision maker has evidence to base outcome on																						
	Probability distributions	distribution of all possible values of random variables and corresponding probabilities																						
	Uncertainty & risk	<ul style="list-style-type: none">A key aspect of solving business problems is dealing with uncertainty EG. Not knowing how much stock to produceProblem solving also involves risk, which depends on the position of the business/decision maker -EG. Interest rates have less effect on businesses that sell essential items to those that sell luxury items																						
	Probability properties & rules	<ul style="list-style-type: none">Event X $0 \leq P(X) \leq 1$Probabilities of all events must = 1 $\sum P(X) = 1$																						
Assigning probabilities	Classical method (priori classical probability)	Occurs with games of chance or where all outcomes are known + probabilities are fixed $\frac{\text{no. of ways events can occur}}{\text{total no. of possible outcomes}}$ <ul style="list-style-type: none">Example: Rolling a die. P(rolling a 4) = 1/6 P(rolling an odd number) = 3/6 P(rolling more than 4) = 2/6 																						
Assigning probabilities	Relative frequency method (empirical classical probability)	Observed (historical) data <ul style="list-style-type: none">past surveys or observations can provide insight to what may occur in the futureit is done through a <u>probability distribution</u> Formula: P(exactly 2) = P(x = 2) = 0.45 (see table →)	<table><caption>Number of insurance claims per day</caption><thead><tr><th>Num of Claims</th><th>Number of Days</th><th>Probability P(X)</th></tr></thead><tbody><tr><td>0</td><td>4</td><td>0.10 = 4/40</td></tr><tr><td>1</td><td>6</td><td>0.15 = 6/40</td></tr><tr><td>2</td><td>18</td><td>0.45 = 18/40</td></tr><tr><td>3</td><td>10</td><td>0.25 = 10/40</td></tr><tr><td>4</td><td>2</td><td>0.05 = 2/40</td></tr><tr><td></td><td>40</td><td>1.00 = 40/40</td></tr></tbody></table> <p>P(Exactly 2) = P(X = 2) = 0.45 P(3 or more in a day) = P(X ≥ 3) = 0.25+0.05 = 0.30</p>	Num of Claims	Number of Days	Probability P(X)	0	4	0.10 = 4/40	1	6	0.15 = 6/40	2	18	0.45 = 18/40	3	10	0.25 = 10/40	4	2	0.05 = 2/40		40	1.00 = 40/40
Num of Claims	Number of Days	Probability P(X)																						
0	4	0.10 = 4/40																						
1	6	0.15 = 6/40																						
2	18	0.45 = 18/40																						
3	10	0.25 = 10/40																						
4	2	0.05 = 2/40																						
	40	1.00 = 40/40																						
Assigning probabilities	Subjective method (guess)	Use this method if there is no method of obtaining probabilities from past experience or mathematical distributions EG. Economic conditions + company's circumstances change rapidly → may be inappropriate to assign probabilities based on historical data																						

		<ul style="list-style-type: none"> may have to adjust probabilities if past figures are no longer suitable in current conditions EG. production manager 'feels' that typically 1 of 500 produced as manufacturing fault Estimate probabilities using data, experience + intuition available This subjective probability value express our degree of belief of what will occur 																																			
Probability rules + properties	Complement of an event A	All outcomes are not part of event A																																			
Probability rules + properties	Joint event (AND)	Involves 2+ characteristics occurring simultaneously																																			
Probability rules + properties	Mutually exclusive events	Events that cannot occur together																																			
Probability rules + properties	Collectively exhaustive events	Set of events such that one of the events must occur																																			
Probability rules + properties	General addition rule (OR)	$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$																																			
Probability rules + properties	Conditional probabilities	$P(A B) = \frac{P(A \wedge B)}{P(B)}$																																			
	If probability of event A is not changed by the existence of event B ...	<p>... then A and B are independent</p> <p>$P(A B) = P(A)$</p>																																			
	Cross-tabs <i>(contingency tables)</i>	<p>sample space for joint events classified by 2 characteristics</p> <ul style="list-style-type: none"> Summarise past data Can be used to provide variety of probability estimates <div data-bbox="510 1241 1120 1548" data-label="Complex-Block"> <p>EXAMPLE: CONROBAR GENDER V PRODUCTIVITY</p> <table border="1"> <thead> <tr> <th></th> <th>Female</th> <th>Male</th> <th>Total</th> </tr> </thead> <tbody> <tr> <td>Prod < 100</td> <td>14</td> <td>16</td> <td>30</td> </tr> <tr> <td>Prod ≥ 100</td> <td>8</td> <td>10</td> <td>18</td> </tr> <tr> <td>Total</td> <td>22</td> <td>26</td> <td>48</td> </tr> </tbody> </table> <ul style="list-style-type: none"> Calculate the following probabilities: • $P(\text{Underperforming employee})$ • $P(\text{Male and underperforming employee})$ • $P(\text{Performing or Female employee})$ • $P(\text{Underperforming given that Male employee})$ • Are productivity performance and gender independent? </div> <div data-bbox="1308 1104 1774 1401" data-label="Table"> <table> <tr> <th rowspan="2">Event</th><th colspan="2">Event</th><th rowspan="2">Total</th></tr> <tr> <th>B₁</th><th>B₂</th></tr> <tr> <td>A₁</td><td>P(A₁ and B₁)</td><td>P(A₁ and B₂)</td><td>P(A₁)</td></tr> <tr> <td>A₂</td><td>P(A₂ and B₁)</td><td>P(A₂ and B₂)</td><td>P(A₂)</td></tr> <tr> <td>Total</td><td>P(B₁)</td><td>P(B₂)</td><td>1</td></tr> </table> <div> <div>Joint Probabilities</div> <div>Simple Probabilities</div> </div> </div>		Female	Male	Total	Prod < 100	14	16	30	Prod ≥ 100	8	10	18	Total	22	26	48	Event	Event		Total	B ₁	B ₂	A ₁	P(A ₁ and B ₁)	P(A ₁ and B ₂)	P(A ₁)	A ₂	P(A ₂ and B ₁)	P(A ₂ and B ₂)	P(A ₂)	Total	P(B₁)	P(B₂)	1	
	Female	Male	Total																																		
Prod < 100	14	16	30																																		
Prod ≥ 100	8	10	18																																		
Total	22	26	48																																		
Event	Event		Total																																		
	B ₁	B ₂																																			
A ₁	P(A ₁ and B ₁)	P(A ₁ and B ₂)	P(A ₁)																																		
A ₂	P(A ₂ and B ₁)	P(A ₂ and B ₂)	P(A ₂)																																		
Total	P(B₁)	P(B₂)	1																																		

	Discrete random variable	<ul style="list-style-type: none">Mean tells us what we expect will happen on averageSD tells us the expected dev around this average																																																																																										
Expected values	Expected value of X (mean):	<div>tells us what we expect will happen on avg</div> <div>formula:</div> <div>$E(x)=\mu=\sum XP(x)$</div> <div><table><tr><td>X</td><td>P(X)</td><td>P(X).X</td></tr><tr><td>0</td><td>0.10</td><td>0.00</td></tr><tr><td>1</td><td>0.15</td><td>0.15</td></tr><tr><td>2</td><td>0.45</td><td>0.90</td></tr><tr><td>3</td><td>0.25</td><td>0.75</td></tr><tr><td>4</td><td>0.05</td><td>0.20</td></tr><tr><td></td><td><u>1.00</u></td><td><u>2.00</u></td></tr></table></div> <div>Expected number of claims/day is: E(X) = 2.00</div>	X	P(X)	P(X).X	0	0.10	0.00	1	0.15	0.15	2	0.45	0.90	3	0.25	0.75	4	0.05	0.20		<u>1.00</u>	<u>2.00</u>																																																																					
X	P(X)	P(X).X																																																																																										
0	0.10	0.00																																																																																										
1	0.15	0.15																																																																																										
2	0.45	0.90																																																																																										
3	0.25	0.75																																																																																										
4	0.05	0.20																																																																																										
	<u>1.00</u>	<u>2.00</u>																																																																																										
Expected values	Expected variation of X (std dev):	<div>tells us the expected variation around the avg</div> <div>$SD(x)=\sigma=\sqrt{(\sum (X-\mu)^2 P(X))}$</div> <div><table><tr><td>X</td><td>X - μ</td><td>(X - μ)²</td><td>P(X)</td><td>(X - μ)²P(X)</td></tr><tr><td>0</td><td>-2</td><td>4</td><td>0.10</td><td>0.40</td></tr><tr><td>1</td><td>-1</td><td>1</td><td>0.15</td><td>0.15</td></tr><tr><td>2</td><td>0</td><td>0</td><td>0.45</td><td>0.00</td></tr><tr><td>3</td><td>1</td><td>1</td><td>0.25</td><td>0.25</td></tr><tr><td>4</td><td>2</td><td>4</td><td>0.05</td><td>0.20</td></tr><tr><td></td><td></td><td></td><td></td><td><u>1.00</u> = σ^2</td></tr></table></div> <div>The SD of claims is $\sqrt{1} = 1$ claim</div>	X	X - μ	(X - μ) ²	P(X)	(X - μ) ² P(X)	0	-2	4	0.10	0.40	1	-1	1	0.15	0.15	2	0	0	0.45	0.00	3	1	1	0.25	0.25	4	2	4	0.05	0.20					<u>1.00</u> = σ^2																																																							
X	X - μ	(X - μ) ²	P(X)	(X - μ) ² P(X)																																																																																								
0	-2	4	0.10	0.40																																																																																								
1	-1	1	0.15	0.15																																																																																								
2	0	0	0.45	0.00																																																																																								
3	1	1	0.25	0.25																																																																																								
4	2	4	0.05	0.20																																																																																								
				<u>1.00</u> = σ^2																																																																																								
Pay off tables	<div>EXAMPLE:</div> <div>Pay off table</div> <div><div><div>Q4.</div><table><tr><th></th><th>Demand</th><th></th><th></th><th></th><th></th></tr><tr><th>Supply</th><th>10</th><th>11</th><th>12</th><th>13</th><th>14</th></tr><tr><td>10</td><td>1000</td><td>1000</td><td>1000</td><td>1000</td><td>1000</td></tr><tr><td>11</td><td>1000-150= 850</td><td>\$1,100</td><td>\$1,100</td><td>\$1,100</td><td>\$1,100</td></tr><tr><td>12</td><td>1000-300= 700</td><td>950</td><td>1200</td><td>1200</td><td>1200</td></tr><tr><td>13</td><td>1000-450= 550</td><td>1100-300= 800</td><td>1200-150= 1050</td><td>1300</td><td>1300</td></tr><tr><td>14</td><td>1000-600= 400</td><td>1100-450= 650</td><td>1200-300= 900</td><td>1300-150= 1150</td><td>1400</td></tr></table></div><div><div>1. Expected Profitability values = (Probability x Profit)</div><table><tr><th>Supply 10'000 eggs</th><th>10</th><th>11</th><th>12</th><th>13</th><th>14</th></tr><tr><td>Profit</td><td>1000</td><td>1000</td><td>1000</td><td>1000</td><td>1000</td></tr><tr><td>Probability</td><td>0.2</td><td>0.2</td><td>0.3</td><td>0.2</td><td>0.1</td></tr><tr><td>Profitability</td><td></td><td></td><td></td><td></td><td></td></tr></table><table><tr><th>Supply 11'000 eggs</th><th>10</th><th>11</th><th>12</th><th>13</th><th>14</th></tr><tr><td>Profit</td><td>850</td><td>1100</td><td>1100</td><td>1100</td><td>1100</td></tr><tr><td>Probability</td><td>0.2</td><td>0.2</td><td>0.3</td><td>0.2</td><td>0.1</td></tr><tr><td></td><td></td><td></td><td></td><td>220</td><td>110</td></tr></table></div></div>		Demand					Supply	10	11	12	13	14	10	1000	1000	1000	1000	1000	11	1000-150= 850	\$1,100	\$1,100	\$1,100	\$1,100	12	1000-300= 700	950	1200	1200	1200	13	1000-450= 550	1100-300= 800	1200-150= 1050	1300	1300	14	1000-600= 400	1100-450= 650	1200-300= 900	1300-150= 1150	1400	Supply 10'000 eggs	10	11	12	13	14	Profit	1000	1000	1000	1000	1000	Probability	0.2	0.2	0.3	0.2	0.1	Profitability						Supply 11'000 eggs	10	11	12	13	14	Profit	850	1100	1100	1100	1100	Probability	0.2	0.2	0.3	0.2	0.1					220	110	
	Demand																																																																																											
Supply	10	11	12	13	14																																																																																							
10	1000	1000	1000	1000	1000																																																																																							
11	1000-150= 850	\$1,100	\$1,100	\$1,100	\$1,100																																																																																							
12	1000-300= 700	950	1200	1200	1200																																																																																							
13	1000-450= 550	1100-300= 800	1200-150= 1050	1300	1300																																																																																							
14	1000-600= 400	1100-450= 650	1200-300= 900	1300-150= 1150	1400																																																																																							
Supply 10'000 eggs	10	11	12	13	14																																																																																							
Profit	1000	1000	1000	1000	1000																																																																																							
Probability	0.2	0.2	0.3	0.2	0.1																																																																																							
Profitability																																																																																												
Supply 11'000 eggs	10	11	12	13	14																																																																																							
Profit	850	1100	1100	1100	1100																																																																																							
Probability	0.2	0.2	0.3	0.2	0.1																																																																																							
				220	110																																																																																							

Supply 12'000 eggs	10	11	12	13	14
Profit	700	950	1200	1200	1200
Probability	0.2	0.2	0.3	0.2	0.1
Profitability x Probability	140	190	360	240	120

Supply 14'000 eggs	10	11	12	13	14
Profit	400	650	900	1150	1400
Probability	0.2	0.2	0.3	0.2	0.1
Profitability x Probability	80	130	270	230	140

Supply 13'000 eggs	10	11	12	13	14
Profit	550	800	1050	1300	1300
Probability	0.2	0.2	0.3	0.2	0.1
Profitability x Probability	110	160	315	260	130

Expected Profit:	
10	\$1,000
11	\$1,050
12	\$1,050
13	\$975
14	\$850

-The manager of the Waverly store is not correct in saying that ordering the maximum possible quantity will return a higher profit.

-This is because the stores who ordered 11000 and 12000 eggs experienced the highest profit at \$1050. Whereas the store who ordered the most eggs (14000) experienced the lowest profit of \$850.

-Probability provides (and data) provides evidence for a decision maker to justify the chosen outcome

Another Example:

Q4 The Colchester Garden Centre purchases and sells Christmas trees during the holiday season. It purchases the trees for \$10 each and sells them for \$20 each. Any trees not sold by Christmas day are sold for \$2 each to a company that makes wood chips. Suppose that the probabilities of the demand for the different number of trees are as follows:

Demand (Number of Trees)	Probability
100	0.20
200	0.50
500	0.20
1000	0.10

Following are payoffs for purchasing 100, 200, 500 or 1000 trees.

a) Complete the missing payoffs;

	Probability	Purchase 100	Purchase 200	Purchase 500	Purchase 1000
Demand 100	0.2	1000	200	-2200	-6200
Demand 200	0.5	1000	2000	-400	-4400
Demand 500	0.2	1000	2000	5000	1000
Demand 1000	0.1	1000	2000	5000	10000

b) Calculate the expected values and suggest how many trees to purchase.

Expected Value

Purchase 100: \$1000

Purchase 200: \$1640

Purchase 500: \$860

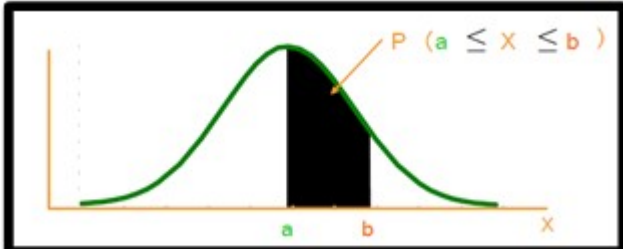
Purchase 1000: -\$2240

<<purchase 200 trees>>

	Mathematical probability distributions	<ul style="list-style-type: none"> • Important source of probabilities • Probabilities can be derived <u>w/o</u> need for data • Each distribution has own unique + important characteristics & properties 	
Type of mathematical probability distribution	Binominal distribution properties	<p><u>4 types of properties:</u></p> <ol style="list-style-type: none"> 1. Fixed no. of identical trials/observations (n) EG. 15 tosses of a coin; 10 light bulbs taken from a warehouse 2. Each trial has only two possible outcomes <ul style="list-style-type: none"> • called 'success' and 'failure' • probability of success is p, probability of failure is q = 1-p EG. head/tail in each toss of a coin; defective/not defective light bulb 3. Constant probability for each trial EG. Probability of a tail is the same each time we toss the coin 4. observations are independent <ul style="list-style-type: none"> • means 1 outcome does not affect another outcome • to ensure independence - observations are randomly sampled 	
Binominal distribution properties	Applications	<ul style="list-style-type: none"> • sampling with replacement • sampling w/o replacement - n < 5% N 	
how to calculate binominal distributions :	Mean	$\mu = E(x) = np$ <p>n = no. of people in sample p = probability</p>	
Binominal distribution characteristics	Variance & Std Dev	$\sigma^2 = np(1-p)$ $\sigma = \sqrt{np(1-p)}$ <p>n = no. of people in sample p = probability</p>	
Binominal formula	Compute probabilities	$P(x) = \frac{n!}{x!(n-x)!} p^x \cdot q^{n-x} \quad \text{for } 0 \leq x \leq n$ <p>Where:</p> <ul style="list-style-type: none"> • n = sample size • p = probability 	

		<ul style="list-style-type: none">$P(x)$ = probability of x success given n and p																																																																																																																																																																																																																																																
Binominal table	Compute probabilities	<p>Demonstration Problem 5.4</p> <p>find binominal probability for:</p> <ul style="list-style-type: none">$n=20$$p=0.4$$x=10$ <table><tr><th rowspan="2">x</th><th colspan="9">$n = 20$</th></tr><tr><th colspan="9">Probability</th></tr><tr><th></th><th>.1</th><th>.2</th><th>.3</th><th>.4</th><th>.5</th><th>.6</th><th>.7</th><th>.8</th><th>.9</th></tr><tr><td>0</td><td>.122</td><td>.012</td><td>.001</td><td>.000</td><td>.000</td><td>.000</td><td>.000</td><td>.000</td><td>.000</td></tr><tr><td>1</td><td>.270</td><td>.058</td><td>.007</td><td>.000</td><td>.000</td><td>.000</td><td>.000</td><td>.000</td><td>.000</td></tr><tr><td>2</td><td>.285</td><td>.137</td><td>.028</td><td>.003</td><td>.000</td><td>.000</td><td>.000</td><td>.000</td><td>.000</td></tr><tr><td>3</td><td>.190</td><td>.205</td><td>.072</td><td>.012</td><td>.001</td><td>.000</td><td>.000</td><td>.000</td><td>.000</td></tr><tr><td>4</td><td>.090</td><td>.218</td><td>.130</td><td>.035</td><td>.005</td><td>.000</td><td>.000</td><td>.000</td><td>.000</td></tr><tr><td>5</td><td>.032</td><td>.175</td><td>.179</td><td>.075</td><td>.015</td><td>.001</td><td>.000</td><td>.000</td><td>.000</td></tr><tr><td>6</td><td>.009</td><td>.109</td><td>.192</td><td>.124</td><td>.037</td><td>.005</td><td>.000</td><td>.000</td><td>.000</td></tr><tr><td>7</td><td>.002</td><td>.055</td><td>.164</td><td>.166</td><td>.074</td><td>.015</td><td>.001</td><td>.000</td><td>.000</td></tr><tr><td>8</td><td>.000</td><td>.022</td><td>.114</td><td>.180</td><td>.120</td><td>.035</td><td>.004</td><td>.000</td><td>.000</td></tr><tr><td>9</td><td>.000</td><td>.007</td><td>.065</td><td>.160</td><td>.160</td><td>.071</td><td>.012</td><td>.000</td><td>.000</td></tr><tr><td>10</td><td>.000</td><td>.002</td><td>.031</td><td>.117</td><td>.176</td><td>.117</td><td>.031</td><td>.002</td><td>.000</td></tr><tr><td>11</td><td>.000</td><td>.000</td><td>.012</td><td>.071</td><td>.160</td><td>.160</td><td>.065</td><td>.007</td><td>.000</td></tr><tr><td>12</td><td>.000</td><td>.000</td><td>.004</td><td>.035</td><td>.120</td><td>.180</td><td>.114</td><td>.022</td><td>.000</td></tr><tr><td>13</td><td>.000</td><td>.000</td><td>.001</td><td>.015</td><td>.074</td><td>.166</td><td>.164</td><td>.055</td><td>.002</td></tr><tr><td>14</td><td>.000</td><td>.000</td><td>.000</td><td>.005</td><td>.037</td><td>.124</td><td>.192</td><td>.109</td><td>.009</td></tr><tr><td>15</td><td>.000</td><td>.000</td><td>.000</td><td>.001</td><td>.015</td><td>.075</td><td>.179</td><td>.175</td><td>.032</td></tr><tr><td>16</td><td>.000</td><td>.000</td><td>.000</td><td>.000</td><td>.005</td><td>.035</td><td>.130</td><td>.218</td><td>.090</td></tr><tr><td>17</td><td>.000</td><td>.000</td><td>.000</td><td>.000</td><td>.001</td><td>.012</td><td>.072</td><td>.205</td><td>.190</td></tr><tr><td>18</td><td>.000</td><td>.000</td><td>.000</td><td>.000</td><td>.000</td><td>.003</td><td>.028</td><td>.137</td><td>.285</td></tr><tr><td>19</td><td>.000</td><td>.000</td><td>.000</td><td>.000</td><td>.000</td><td>.000</td><td>.007</td><td>.058</td><td>.270</td></tr><tr><td>20</td><td>.000</td><td>.000</td><td>.000</td><td>.000</td><td>.000</td><td>.000</td><td>.001</td><td>.012</td><td>.122</td></tr></table> <p>$n = 20$</p> <p>$p = 0.40$</p> <p>$P(X = 10) = 0.117$</p>	x	$n = 20$									Probability										.1	.2	.3	.4	.5	.6	.7	.8	.9	0	.122	.012	.001	.000	.000	.000	.000	.000	.000	1	.270	.058	.007	.000	.000	.000	.000	.000	.000	2	.285	.137	.028	.003	.000	.000	.000	.000	.000	3	.190	.205	.072	.012	.001	.000	.000	.000	.000	4	.090	.218	.130	.035	.005	.000	.000	.000	.000	5	.032	.175	.179	.075	.015	.001	.000	.000	.000	6	.009	.109	.192	.124	.037	.005	.000	.000	.000	7	.002	.055	.164	.166	.074	.015	.001	.000	.000	8	.000	.022	.114	.180	.120	.035	.004	.000	.000	9	.000	.007	.065	.160	.160	.071	.012	.000	.000	10	.000	.002	.031	.117	.176	.117	.031	.002	.000	11	.000	.000	.012	.071	.160	.160	.065	.007	.000	12	.000	.000	.004	.035	.120	.180	.114	.022	.000	13	.000	.000	.001	.015	.074	.166	.164	.055	.002	14	.000	.000	.000	.005	.037	.124	.192	.109	.009	15	.000	.000	.000	.001	.015	.075	.179	.175	.032	16	.000	.000	.000	.000	.005	.035	.130	.218	.090	17	.000	.000	.000	.000	.001	.012	.072	.205	.190	18	.000	.000	.000	.000	.000	.003	.028	.137	.285	19	.000	.000	.000	.000	.000	.000	.007	.058	.270	20	.000	.000	.000	.000	.000	.000	.001	.012	.122	
x	$n = 20$																																																																																																																																																																																																																																																	
	Probability																																																																																																																																																																																																																																																	
	.1	.2	.3	.4	.5	.6	.7	.8	.9																																																																																																																																																																																																																																									
0	.122	.012	.001	.000	.000	.000	.000	.000	.000																																																																																																																																																																																																																																									
1	.270	.058	.007	.000	.000	.000	.000	.000	.000																																																																																																																																																																																																																																									
2	.285	.137	.028	.003	.000	.000	.000	.000	.000																																																																																																																																																																																																																																									
3	.190	.205	.072	.012	.001	.000	.000	.000	.000																																																																																																																																																																																																																																									
4	.090	.218	.130	.035	.005	.000	.000	.000	.000																																																																																																																																																																																																																																									
5	.032	.175	.179	.075	.015	.001	.000	.000	.000																																																																																																																																																																																																																																									
6	.009	.109	.192	.124	.037	.005	.000	.000	.000																																																																																																																																																																																																																																									
7	.002	.055	.164	.166	.074	.015	.001	.000	.000																																																																																																																																																																																																																																									
8	.000	.022	.114	.180	.120	.035	.004	.000	.000																																																																																																																																																																																																																																									
9	.000	.007	.065	.160	.160	.071	.012	.000	.000																																																																																																																																																																																																																																									
10	.000	.002	.031	.117	.176	.117	.031	.002	.000																																																																																																																																																																																																																																									
11	.000	.000	.012	.071	.160	.160	.065	.007	.000																																																																																																																																																																																																																																									
12	.000	.000	.004	.035	.120	.180	.114	.022	.000																																																																																																																																																																																																																																									
13	.000	.000	.001	.015	.074	.166	.164	.055	.002																																																																																																																																																																																																																																									
14	.000	.000	.000	.005	.037	.124	.192	.109	.009																																																																																																																																																																																																																																									
15	.000	.000	.000	.001	.015	.075	.179	.175	.032																																																																																																																																																																																																																																									
16	.000	.000	.000	.000	.005	.035	.130	.218	.090																																																																																																																																																																																																																																									
17	.000	.000	.000	.000	.001	.012	.072	.205	.190																																																																																																																																																																																																																																									
18	.000	.000	.000	.000	.000	.003	.028	.137	.285																																																																																																																																																																																																																																									
19	.000	.000	.000	.000	.000	.000	.007	.058	.270																																																																																																																																																																																																																																									
20	.000	.000	.000	.000	.000	.000	.001	.012	.122																																																																																																																																																																																																																																									
Binominal distribution	EXAMPLE: binominal distribution	<p>A components manufacturer samples 20 products randomly from each shift to monitor quality. Over an extended period they have found that the number of defectives is approximately 10% of all products made.</p> <p>Explain why the Binomial distribution is appropriate in this case.</p> <p>$n = 20$ number of trials (products)</p> <p>Two possible outcomes for each trial (defective/not defective)</p> <p>$P(\text{defective}) = 0.10$; $P(\text{not defective}) = 0.90$</p> <p>Products independent (randomly selected)</p> <p>What would be the expected number of defectives in each shift.</p> <p>$E(X) = n \times p = 20 \times 0.10 = 2$ products</p> <p>In a particular shift, five (5) defectives were found. Find the probability of this occurring (or worse). Comment on your answer.</p> <p>Want $P(X \geq 5) = 1 - P(X < 5)$</p> <p>$= 1 - 0.957$</p> <p>$= 0.043$</p> <p>Small chance (4.3%) to get five or more defective products.</p> <p>There may be a problem with this particular shift.</p>																																																																																																																																																																																																																																																

how to recognise Poisson distributions	Poisson distribution properties	<p><u>4 essential properties:</u></p> <ol style="list-style-type: none"> You wish to count no. of occurrences over an interval: <ul style="list-style-type: none"> Each occurrence is independent of other occurrences No. of occurrences in each interval can range from 0 to ∞ Average (expected) no. of occurrences in interval is λ (lambda) <p>EG.</p> <ul style="list-style-type: none"> seeing how many pixel burnouts there are in TV surface area No. of telephone calls per minute at a small business No. of customer arrivals at a bank in an hour No. of paint spots per new vehicle No. of units of product demanded per week probability that an event occurs in 1 area of opportunity = same for all areas of opportunity EG. Pixel burnout probability for TV screen is same for middle of screen, top + bottom no. of events that occur in 1 area of opportunity independent of no. of events that occur in other areas of opp EG. 1 computer crash will not affect another computer crash probability that 2+ events occur in an area of opportunity approaches zero as the area of opportunity becomes smaller EG. when you focus whole TV there is greater probability of pixel burnout - if focus only on bottom of TV = less likely to get burnout 	
Calculate poisson distributions	Poisson distribution characteristics	<p>Mean: $\mu = \lambda$</p> <p>Variable + Std Dev: $\sigma^2 = \lambda$ $\sigma = \sqrt{\lambda}$</p> <p>Where λ = expected no. of occurrences in the interval</p>	
Poisson formula	Calculate probabilities	<p>$P(x) = \frac{e^{-\lambda} \lambda^x}{x!}$ For $x = 0, 1, 2, 3, \dots$</p> <p>Where: λ = long – run average $E = 2.718282 \dots$ (base of natural logarithms)</p>	

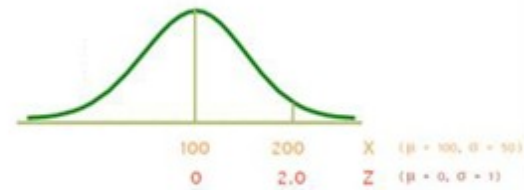
Poisson distribution	EXAMPLE: Poisson distribution	<p>The number of times a company's HR system "crashes" is, on average, 3.3 per month.</p> <p>Explain why the Poisson distribution would apply in this case.</p> <p>Interested in number of times a crash occurs over an interval (one month period).</p> <p>Assume that a crash in a one month interval has no effect on any other crash during any one month interval (independent).</p> <p>Number of crashes in a month has no upper limit.</p> <p>Expected number of crashes (lambda) per month is known.</p> <p>Calculate the probability that the system does not crash in a given month.</p> <p>Want $P(X = 0) = 0.0369$</p> <p>What is the probability it does crash three or more times in a given month.</p> <p>Want $P(X \geq 3) = P(X=3) + P(X=4) + P(X=5) + \dots$ $= 1 - P(X < 3)$ $= 1 - 0.3594 = 0.6406$</p>	
	Probability & decision making	<ul style="list-style-type: none"> Probabilities can be used to measure lvl of uncertainty involved Use of probabilities in decision making does not necessarily mean we get a desired outcome Common mistake is to associate good (poor) decisions with good (poor) outcomes May be the case that a persons a sound decision (w high probability) but due to uncertainty of situation = unlucky + diff outcome subsequently occurred 	
Normal distributions	characteristics of normal distributions	<ul style="list-style-type: none"> bell shaped symmetrical mean, median + mode are equal location is determined by mean spread is determined by std dev depends on 2 parameters: Mean + std dev 	
Normal distributions	calculate normal distribution	<ul style="list-style-type: none"> probability measured by area under curve: 	

-Need to transform X units into Z units:

$$Z = \frac{X - \mu}{\sigma}$$

The Z distribution has mean = 0 and standard deviation = 1

COMPARING X AND Z UNITS



Note that the distribution is the same, only the scale has changed. We can express the problem in original units (X) or in standardised units (Z)

EXAMPLE:
normal
distribution

If X is distributed normally with mean of 100 and standard deviation of 50, then:

$$Z = \frac{X - \mu}{\sigma} = \frac{200 - 100}{50} = 2.0$$

This says that X = 200 is two standard deviations above the mean of 100

GENERAL PROCEDURE FOR FINDING PROBABILITIES

To find $P(a < X < b)$ when X is distributed normally:



1. Draw the normal curve for the problem in terms of X
2. Translate X-values to Z-values
3. Use the Standardised Normal Table or software to find the area

FINDING THE X VALUE FOR A KNOWN PROBABILITY

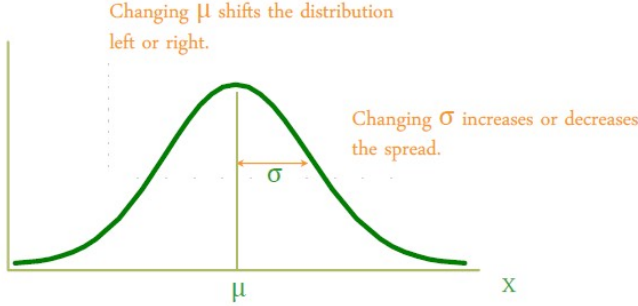
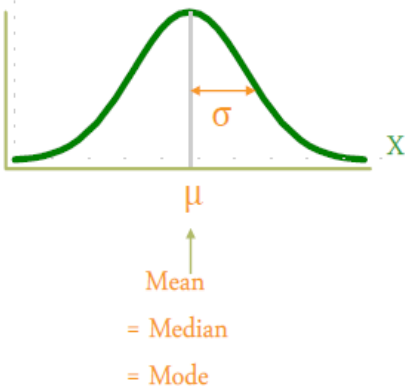
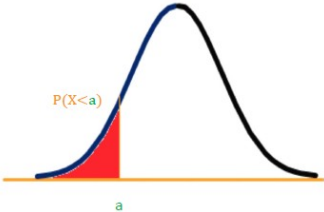
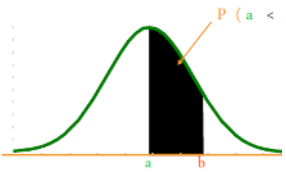
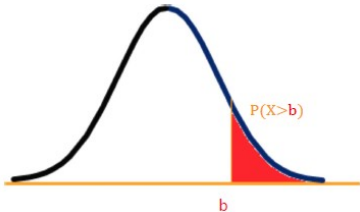
Steps to find the X value for a known probability:

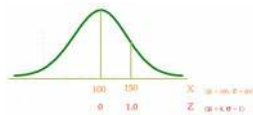
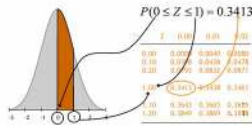
1. Draw the normal curve for the problem showing the known probability
2. Find the Z value for the known probability
3. Convert to X units using the formula


$$X = \mu + Z\sigma$$


	<p>EXAMPLE: finding probabilities</p>	<div data-bbox="557 89 1211 403" data-label="Complex-Block"> <p>Weekly demand for a particular brand of a perishable product is assumed to be approximately normal, with a mean of 150 and a standard deviation of 8. Management orders the product at the beginning of the week and throws out leftover stock at the end of the week.</p> <p>What would you recommend as the minimum and maximum order quantities for this product? Explain.</p> <p>Minimum = $\mu - 3\sigma = 150 - 3 \times 8 = 126$</p> <p>Maximum = $\mu + 3\sigma = 150 + 3 \times 8 = 174$</p>  </div>	
	<p>EXAMPLE: finding X value for known probability</p>	<div data-bbox="551 737 1256 1058" data-label="Complex-Block"> <p>What should the ordering policy be if they want to ensure that they run out of stock ('stockout') in only about 5% of weeks?</p> <p>Want $P(X > ?) = 0.05$</p> <p>$Z = 1.645$ (from table)</p> <p>$Z = (X - \mu) / \sigma = (X - 150) / 8 = 1.645$</p> <p>$\therefore X - 150 = 8 \times 1.645 = 13.16$</p> <p>$\therefore X = 13.16 + 150 = 163.16$</p> <p>Company should order at least 164 products per week</p>  </div>	

Week 6 (Topic 6) - Chapter FIVE (Continuous Probability Distribution)

	Continuous distributions	<ul style="list-style-type: none"> associated with random variables that take values at any point over given interval uniform + exponential distributions are important because normal distribution is most widely encountered continuous distribution 	
	continuous probability distributions	<ul style="list-style-type: none"> no longer talk of probability of random variable being equal to specific value talk about probability of random variable having value within given range area under probability curve b/w 2 given values = the probability that random variable lies b/w the 2 values 	
	characteristics of normal distribution	<ol style="list-style-type: none"> Bell-shaped Symmetrical distribution Continuous distribution Mean = median = mode (all equal) Location (centre) is determined by the mean (μ) Spread is determined by the standard deviation (σ) Area under the curve is equal to 1 <div style="display: flex; align-items: center;">  <div style="margin-left: 20px;">  </div> </div>	
finding normal probabilities	Normal probability	<ol style="list-style-type: none"> measured by area under the curve <div style="display: flex; justify-content: space-around; align-items: center;">    </div>	
	Standardised normal distribution	<ol style="list-style-type: none"> any normal distribution (with any mean + std dev combination) can be transformed into standardised normal distribution (Z) need to transfer X units into Z units 	

		<p>Formula:</p> $Z = \frac{X - \mu}{\sigma}$ <p>* standardised normal (z) distribution has mean = 0 std dev = 1</p>	
Standardise d normal distribution	EXAMPLE: stdardised normal distribution	<p>2. If X is distributed normally with mean of 100 and std dev of 50, z-score for X = 150 is:</p> $z = \frac{X - \mu}{\sigma} = \frac{150 - 100}{50} = 1.0$ <p>3. Z-score = no. of std dev particular value (X) is away from mean This says that X = 150 is 1 std dev above mean of 100</p>	
	Comparing X & Z	<p>distribution is the same</p> <ul style="list-style-type: none"> scale has changed express problem in standardised units (Z) 	
Procedure for finding probabilities	Find P(a < X < b) when X is distributed normally	<p>Find P(a < X < b) when X is distributed normally:</p> <ol style="list-style-type: none"> draw normal curve for problem in terms of X translate X-values to Z-values Use standardised Normal Table or software to find the area (probability) 	
	Table look up of std normal probability		
Finding X for known probability (Reverse	Finding X for known probability (Reverse	<p>Steps to find X value for known probability:</p> <ol style="list-style-type: none"> draw the normal curve for probalm showing area of interest find probability (area) among values in body of table determine value of Z (in left tail Z must have (neg) -sign, right tail Z has (pos) +sign) 	

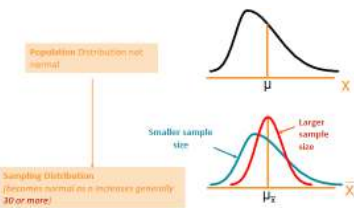
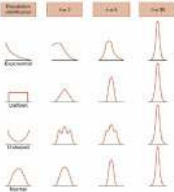
Questions)	Question)	4. calculate $X = Z\sigma + \mu$ Calculate $X = Z\sigma + \mu$	
Rules	Reading the table	<ul style="list-style-type: none"> If exact probability is not in body of table - find closest probability + determine Z value to 2 decimal places <u>DO NOT</u> interpolate except when value is exactly HALF way EG. 0.45 is exactly halfway between 0.4495 (for Z=-1.64) & 0.4505 (for Z=-1.65) hence use Z= -1.645 	
	EXAMPLE	<ul style="list-style-type: none"> Weekly demand for a particular brand of perishable product = normal Mean = 150 Std Dev = 8. <p>Management orders the product at the beginning of the week and throws out leftover stock at the end of the week. What would you recommend as the min + max order quantities for this product? Explain.</p> <p>Min: $\mu - 3\sigma$ $150 - 3 \times 8$ $\hookrightarrow 126$</p> <p>Max: $\mu + 3\sigma$ $150 + 3 \times 8$ $\hookrightarrow 174$</p> <p>The ordering policy has been to order 160 units of the brand. What is the probability that some product be thrown out at the end of the week?</p> <p> $P(\text{Wastage}) = P(\text{Demand} < 160)$  $z = \frac{X - \mu}{\sigma}$ $\frac{160 - 150}{8}$ $\frac{10}{8}$ $\hookrightarrow 1.25$ $P(z < 1.25) = 0.8944$ (from table) </p> <p>= there is almost 90% chance that there will be at least some wastage</p>	


		<p>What should ordering policy be if they want to ensure that they run out of stock in only 5% of week?</p>  <p>Want: $P(X > ?) = 0.05$ $Z = 1.645$ (table) $Z = \frac{X - \mu}{\sigma}$ $\hookrightarrow \frac{X - 150}{8}$ $\hookrightarrow 1.645$ $\therefore X - 150$ $\hookrightarrow 8 \times 1.645$ $\hookrightarrow 13.16$ $\therefore X = 13.16 + 150$ $\hookrightarrow 163.16$</p>	
Summary		<ul style="list-style-type: none"> • normal distribution most important distribution in statistics due to wide application in solving many practical & business related problems • probabilities of continuous random variables - calculated as areas under curve & b/w specified values of random variable • normal distribution with mean μ & std dev σ can be transformed to std normal distribution involving z-scores • use z-scores allow probabilities to calculate from standardised normal distribution table 	

Week 7 (Topic 7) - Chapter SEVEN (Sampling & Sampling Distributions)			
5 key words in Statistics	5 key words in Statistics	<ol style="list-style-type: none"> 1. Population 2. Census 3. Sample 4. Error 5. Probability <p>**Population is the “goal” – our purpose – even though we spend most of the time working with the sample.</p>	
	Types of Samples used	<pre> graph TD Samples --> NonRandom[Non-Random Samples] Samples --> Random[Random Samples] NonRandom --> SelfSelection[Self-selection] NonRandom --> Convenience[Convenience] NonRandom --> Judgment[Judgment] Random --> SimpleRandom[Simple Random] Random --> Systematic[Systematic] Random --> Stratified[Stratified] Random --> Cluster[Cluster] </pre>	
	Non-Random (Non-Probability) Sampling	<p>Probability of particular element of population entering sample is unknown</p> <ul style="list-style-type: none"> • some individuals / items in population have greater chance of selection than others <p>= statistically valid statements or inferences <u>cannot</u> be made about precision of estimates</p>	
	Advantages of non-random sampling:	<ul style="list-style-type: none"> • sampling costs are lower + implementation is easier • sometimes there is no alternative (random sample cannot be taken) <p>However this means:</p> <ul style="list-style-type: none"> • inferential statistical techniques <u>cannot</u> be applied • no statistical statement should be made re precision of result 	
	Self-selection	people are invited to submit questionnaire EG. online	
	Convenience	elements included in sample are chosen b/c of accessibility or willingness EG. students are interviewed as they enter library	
	Judgment	<ul style="list-style-type: none"> • knowledgeable person selects sampling units that he/she feels are most representative of population • quality of result dependent on judgement of person selecting sample • person may be well-qualified - judgement may be highly regarded = does not give results statistical validity 	
	random	<ul style="list-style-type: none"> • each element of population has known (non-zero) chance of being included in sample chosen 	

	NB: assume all samples are random for this unit	<ul style="list-style-type: none"> • should be used where possible • inferential statistics requires random samples. 	
	simple random sampling	<ul style="list-style-type: none"> • every individual/item from sampling frame (list of eligible elements from population) has equal chance of being selected • every sample of fixed size has same chance of selection as every other sample of that size • most elementary random sampling technique • samples normally obtained from computer random number generators 	
	Systematic sampling	<ul style="list-style-type: none"> • decide on sample size • divide frame of N individuals into groups of K individuals: $k = N/n$ • Randomly select 1 individual from 1st group • select every k^{th} individual thereafter • some situations it is more convenient/faster than simple random sample 	
	stratified sampling	<ul style="list-style-type: none"> • divide population in 2+ subgroups ("strata") according to common characteristics EG. gender, age groups, states etc • simple random sample is selected from each sub group ("stratum") with sample sizes proportional to strata sizes • samples from each stratum are combined into one 	
	cluster sampling	<ul style="list-style-type: none"> • population is divided into several "clusters" each rep of population • simple random sample of clusters is selected - all items in cluster can be used or items can be chosen from cluster using another probability sampling techniques • clusters are often based on geographical groupings where cost of accessing individual households can be reduced EG. postcodes, electorates, city blocks 	
survey errors	non- sampling error	<p>range of errors either for census or sample include:</p> <ul style="list-style-type: none"> • response + non-response bias • interview bias • self-selection bias • measurement error • coverage error • processing error (typos) 	
	sampling error	<ul style="list-style-type: none"> • almost certain to involve error • a (relatively small) sample is unlikely to have exactly same features as population which it was drawn • RECALL: <ul style="list-style-type: none"> - parameter is numerical characteristic of pop - statistic is numerical characteristic of sample 	

	point estimation	<ul style="list-style-type: none"> from sample, value of relevant sample statistic could be used as point estimates of equivalent pop parameters \hat{x} = point estimator of population mean (μ) s = point estimator of population std dev (σ) \hat{p} = point estimator of population proportion (p) 	
	sampling error	<ul style="list-style-type: none"> diff b/w sample statistic estimate + corresponding population parameter sampling errors are: $\hat{x} - \mu$ = sample mean $s - \sigma$ = sample std dev $\hat{p} - p$ = sample proportion manage error using probability 	
	Manage estimate sample error	<ul style="list-style-type: none"> using random sampling combined with probability = control + predict magnitude of error likely to occur statistical theory based on concept known as sampling distribution of sample stat show relatively small samples (at great savings of time + money) can provide remarkably high degrees of accuracy in estimating features of population 	
	sampling distribution	<ul style="list-style-type: none"> distribution of all possible values of statistic for given size sample selected from population 	
	Central limit theorem	<div data-bbox="510 762 797 916" data-label="Figure"> </div> <ul style="list-style-type: none"> if population distribution is normal = sample distribution of \hat{x} will be normal <u>regardless of</u> size of n if population distribution is <u>not</u> normal = sample distribution of \hat{x} will be normal or approximately normal when n is sufficiently large enough (30+) <p>in both cases: $\mu_{\hat{x}} = \mu$</p> <div data-bbox="510 1382 658 1538" data-label="Diagram"> </div> <div data-bbox="716 1369 927 1481" data-label="Figure"> </div>	
	EXAMPLE:		

	Normal population		
	EXAMPLE: Not normal population		
	Distribution of sample means from various sample sizes		
	Z-formula for	<ul style="list-style-type: none"> z-score is used to calculate probabilities from sampling distribution of \bar{x} 	

	sampling distribution of \bar{X}	<ul style="list-style-type: none"> z-formula for the sampling distribution of \bar{X} $Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$ $\hat{=} \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$ <p>Where: \bar{X} = sample mean μ = population mean σ = population std dev n = sample size</p>	
	Std error of sample mean	<ul style="list-style-type: none"> measure of variability in the mean from sample to sample given by std error of the sample mean (std dev of all possible sample mean). NB: std error of sample mean <u>decreases</u> as the sample size n <u>increases</u> Std error = avg error expected to make in using sample mean as point estimate of population mean <p>Formula: $\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$</p>	
	EXAMPLE	<ul style="list-style-type: none"> A population has mean $\mu = 8$ Std Dev $\sigma = 3$ A random sample of size $n = 36$ is selected. <p>What is the probability that the sample mean is between 7.8 and 8.2?</p> <p>Even if the population is not normally distributed, the central limit theorem can be used as $n > 30$</p> <ul style="list-style-type: none"> sampling distribution of \bar{X} = approximately normal mean $\mu_{\bar{X}} = 8$ std dev:  $\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$ $\hat{=} \frac{3}{\sqrt{36}}$ $\hat{=} 0.5(\text{std error})$	

	Sampling distribution of \hat{p}	<ul style="list-style-type: none"> is the probability distribution of all possible values of the sample proportion \hat{p} for given sample size (n) <p>sample proportion: $\hat{p} = \frac{x}{n}$</p> <p>where: x = no. of items in sample with characteristics n = no. of items in sample</p> <p>sample distribution:</p> <ul style="list-style-type: none"> approx. normal = $np > 5$ & $nq > 5$ p = population proportion & $q = 1 - p$ std dev of distribution = $\sqrt{\frac{pq}{n}}$ 	
	Std error of sample proportion (\hat{p})	<ul style="list-style-type: none"> measure of variability in proportion from sample to sample given by std error of sample proportion $\sigma_{\hat{p}} = \sqrt{\frac{pq}{n}}$ <ul style="list-style-type: none"> NB: std error of sample proportion decreases as sample size (n) increases Std error is the avg error we expect to make using sample proportion as point estimate of population proportion 	
	Z-formular for sampling distribution of \hat{p}	<ul style="list-style-type: none"> When $np > 5$ & $nq > 5$ $z = \frac{\hat{p} - p}{SE_{\hat{p}}}$ $= \frac{\hat{p} - p}{\sqrt{\frac{pq}{n}}}$ <p>Where: \hat{p} = sample proportion n = sample size p = population proportion $q = 1 - p$</p>	
	EXAMPLE	<ul style="list-style-type: none"> If true proportion of voter who support Proposition A is $p = 0.4$ What is possibility that sample size of 200 yields sample proportion b/w 0.40 & 0.45? If $p = 0.4$ & $n = 200$ what is: $P(0.40 \leq \hat{p} \leq 0.45)$? 	

Find $\sigma_{\hat{p}}$: $\sigma_{\hat{p}} = \sqrt{\frac{pq}{n}} = \sqrt{\frac{0.4(1-0.4)}{200}} = 0.03464$

Convert to standardised normal:



$$P(0.40 \leq \hat{p} \leq 0.45) = P\left(\frac{0.40 - 0.40}{0.03464} \leq Z \leq \frac{0.45 - 0.40}{0.03464}\right) = P(0 \leq Z \leq 1.44)$$

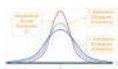

Use standardised normal table: $P(0 \leq Z \leq 1.44) = 0.4251$





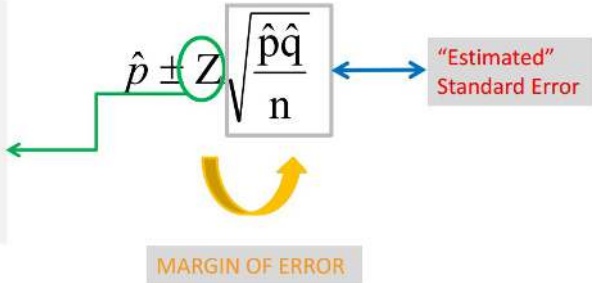
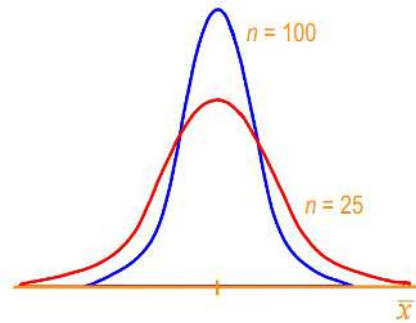
Week 8 (Topic 8) - Chapter EIGHT (Confidence Intervals)

sampling distributions	<p>EXAMPLE male v female productivity</p> <ul style="list-style-type: none"> sample mean for productivity: <ul style="list-style-type: none"> males: $\hat{x}_M = 99.9\%$ females: $\hat{x}_F = 97.5\%$ are 26 males & 22 females enough to draw conclusions re all 3,000 employees at Conrobar? what is mean (avg) productivity: <ul style="list-style-type: none"> for <u>all</u> males? $\mu_M = ?$ for <u>all</u> females? $\mu_F = ?$ 	
Point estimation	<ul style="list-style-type: none"> \hat{x} = point estimator of population mean (μ) s = point estimator of population std dev (σ) \hat{p} = point estimator of population proportion (p) 	
Why don't use point estimate	<ul style="list-style-type: none"> Should <u>not</u> use value of \hat{x} by itself as estimate of μ as: <ul style="list-style-type: none"> Almost certain to be wrong, and We won't know how wrong. Point estimator <u>does</u> not provide info re how close estimate is to population parameter We would have <u>no confidence</u> in using it – instead we use confidence interval 	
Margin of error	<ul style="list-style-type: none"> Interval estimate is constructed by subtracting + adding margin of error (ME) to point estimate: <i>Sample statistic \pm ME</i> 	
Interval estimate	<ul style="list-style-type: none"> Interval estimate of population mean is: $\hat{x} \pm ME$ 	
	<p>Example of 95% interval estimate</p>	

	<p>Confidence interval of μ (σ known) <i>i</i></p>	<ul style="list-style-type: none"> • Use z (stdarised normal) distribution • General form of our confidence interval would be: <ul style="list-style-type: none"> - use the z (stdised normal) distribution - general form of confidence interval would be:  <ul style="list-style-type: none"> • Most commonly used confidence intervals + corresponding z values: <ul style="list-style-type: none"> - 90% of values within ± 1.645 std error of μ - 95% of values within ± 1.96 std error of μ - 98% of values within ± 2.33 std error of μ - 99% of values within ± 2.575 std error of μ • if you construct a 95% confidence interval, 95% is called confidence coefficient other intervals can be used (eg. 96%) 													
	<p>EXAMPLE: no. of defects</p>	<ul style="list-style-type: none"> • calculate 95% confidence interval estimate for mean no. of defects per shift across <u>all</u> shifts • assume population & std dev is known: <ul style="list-style-type: none"> - $\sigma = 10$ - $\bar{x} = 11.6$ - $n = 40$ - 95% confidence = 1.96 std errors (z)  <p>Interpretation: we are 95% confident that the mean no. of defects per shift across all shifts is somewhere in the range 8.5 to 14.7 defects.</p> <ul style="list-style-type: none"> • calculate 90% & 99% confidence interval estimate for true mean no. of defects per shift <table border="1" data-bbox="607 1315 1216 1382"> <thead> <tr> <th colspan="2">90% Interval: μ</th> <th colspan="2">99% Interval: μ</th> </tr> </thead> <tbody> <tr> <td>- From:</td> <td>8.999</td> <td>- From:</td> <td>7.527</td> </tr> <tr> <td>- To:</td> <td>14.201</td> <td>- To:</td> <td>15.673</td> </tr> </tbody> </table> <ul style="list-style-type: none"> • comparing intervals to 95% CI, can see that there is trade-off b/w confidence & margin of error **higher the confidence = less precise (or wider) the interval is 	90% Interval: μ		99% Interval: μ		- From:	8.999	- From:	7.527	- To:	14.201	- To:	15.673	
90% Interval: μ		99% Interval: μ													
- From:	8.999	- From:	7.527												
- To:	14.201	- To:	15.673												

	T Distribution	$t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$ $z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$ <ul style="list-style-type: none"> Shape of distribution similar to normal distribution <u>but</u> varies re sample size 	
T-distribution	Properties	<ul style="list-style-type: none"> T-scores are bigger than equivalent z-scores Family of distributions T-score depends on degree of freedom (d.f.) Bell shaped like normal distribution – <u>but</u> flatter + wider T-distribution provides more conservative confidence interval estimate, since intervals (slightly) wider for given sample size <p>NB: T-distribution should be used if σ is unknown T \rightarrow Z as n increases</p> 	
	Confidence interval for μ (σ unknown)	<ul style="list-style-type: none"> Use T-distribution General form of confidence interval would be: <ul style="list-style-type: none"> use t-distribution general form of our confidence interval would be: 	
	EXAMPLE:	<ul style="list-style-type: none"> Assume population std dev (σ) is unknown (more realistic) 	

	No. of Defects	<ul style="list-style-type: none"> • Use value of sample std dev (s) (9.77) as point estimate of σ • Use T-distribution instead of normal distribution • Calculate 95% confidence interval estimate for true mean no. of defects per shift 	
Conrobar Productivity	EXAMPLE: Confidence Interval	<ul style="list-style-type: none"> • 95% confidence intervals for the mean productivity per group (male/female) as follows: • All males: <ul style="list-style-type: none"> - Sample mean: $\bar{x} = 99.981\%$ - 95% confidence interval: 98.376% to 101.585% • All females: <ul style="list-style-type: none"> - Sample mean: $\bar{x} = 97.53\%$ - 95% confidence interval: 95.499% to 99.574% • Males: 95% confidence interval: 98.376% to 101.585% All males could be achieving 100% on average • Females: 95% confidence interval: 95.499% to 99.574% All females not achieving 100% on avg. • At 95% confidence, we <u>cannot</u> conclude males are doing better on avg than females • 98% or 99% confidence intervals are wider: same conclusions except females could now be achieving 100% avg • 90% confidence intervals are narrower: different conclusions 	

<p>Categorical variables</p> <p>(assume ME either side of sample proportion)</p>	<p>Confidence interval for p</p>	<p>general form of interval is:</p> <div data-bbox="573 129 1406 413"> <p>z indicates how <u>wide</u> the confidence interval is in terms of the <u>number of standard errors</u>. Eg. 1.96 std. errors wide. The z value ties back to the level of confidence.</p>  <p>MARGIN OF ERROR</p> </div>	
	<p>EXAMPLE: Shifts with defects</p>	<p>What is true proportion of shifts which defects occurred (use 95% confidence)</p> <ul style="list-style-type: none"> Of sample of 40 shifts – defects found across 33 shifts: $\hat{p} = \frac{33}{40} = 0.825$ Sample size lrg enough for sampling distribution of \hat{p} to be approximated by normal distribution 95% confidence = 1.96 std errors (z) 95% confidence interval: $\hat{p} \pm z \sqrt{\frac{\hat{p}\hat{q}}{n}} = 0.825 \pm 1.96 \sqrt{\frac{0.825 \times 0.175}{40}}$ $= 0.825 \pm 1.96 \times 0.06$ $= 0.825 \pm 0.118$ $= 0.707 \text{ to } 0.943$ 95% confident true proportion of shifts where there are defects somewhere in range: 70.7% to 94.3% 	
	<p>Different sample sizes</p>	<ul style="list-style-type: none"> Logic tells us the bigger the sample = more accurate If n is larger \therefore SE must be smaller: $SE = \frac{s}{\sqrt{n}}$ Sample distribution would be narrower Confidence interval \therefore be narrower <div data-bbox="1339 1145 1753 1469">  </div>	

Calculating sample size (n)	<ul style="list-style-type: none"> Before taking sample to estimate μ or p – first estimate what size sample required To do this, specify: <ul style="list-style-type: none"> Margin of error (ME) Degree of confidence required (eg. 95%) Obtain estimate of σ and p/q <div style="display: flex; justify-content: space-around; align-items: center;"> <div style="text-align: center;"> $n = \frac{z^2 \sigma^2}{ME^2}$ <div style="background-color: #f0f0f0; padding: 2px; border: 1px solid #ccc;">For numerical data</div> </div> <div style="text-align: center;"> $n = \frac{z^2 pq}{ME^2}$ <div style="background-color: #f0f0f0; padding: 2px; border: 1px solid #ccc;">For categorical data</div> </div> </div>	
EXAMPLE: No. of defects	<ul style="list-style-type: none"> If we know mean no. of defects per shift (3 defects) + 95% confidence Use $\sigma = 10$ $n = \frac{z^2 \sigma^2}{ME^2} = \frac{1.96^2 10^2}{3^2} = 42.7$ $N = 43$: is min sample size needed for specifications 	NB: always round up
EXAMPLE: Shifts with defects	<ul style="list-style-type: none"> if want to know true proportion of shifts in which defects occur to within 5% & 95% confidence use previous study's $\hat{p} = 0.825$ to approx. p $n = \frac{z^2 pq}{ME^2} = \frac{1.96^2 (0.825)(0.175)}{0.05^2} = 221.84$ <p>NB: $n = 222$ is min sample size needed for our specifications</p>	NB: always round up
Calculating sample size	<ul style="list-style-type: none"> if no idea re likely value of p then use $p = 50\%$ in sample size calculation if some idea re likely value of p then: <ul style="list-style-type: none"> pilot study similar / previous study <p>\therefore use that in calculation of n</p>	
EXAMPLE: Newspaper polls	<ul style="list-style-type: none"> election poll is estimate with 95% confidence, proportion of voters who will vote for labor party within 3% accuracy use $p = 0.5$ (no better estimate of p is available) $n = \frac{z^2 pq}{ME^2} = \frac{1.96^2 (0.5)(0.5)}{0.03^2} = 1067.1$ <p>NB: $n = 1068$ is min sample size needed for specifications</p>	

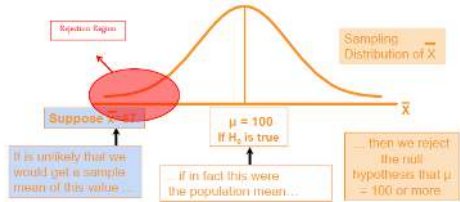
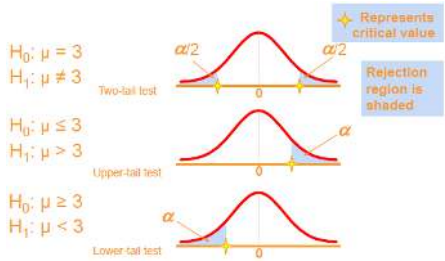
Week 9 (Topic 9) - Chapter NINE (Hypothesis Testing)			
	EXAMPLE: Conrobar	<ul style="list-style-type: none"> Sample of 48 employees showed: <ul style="list-style-type: none"> Mean productivity = 98.9% Mean days absent = 2.4 days Do these mean ALL employees, on average, are not meeting the 100% standard for productivity and the 1.5 days target for days absent? Need to allow for sampling error Confidence intervals (last lecture) provide one approach Hypothesis tests provide another 	
2 key inferential tools	Confidence interval estimation	use if no idea re value of population parameter being investigated	
2 key inferential tools	Hypothesis test	Use if some idea re value of population parameter being investigated or if some hypothesised value against can compare sample results	
2 hypotheses	Null Hypothesis (H_0)	assumed population parameter is <u>correct</u>	
	Alternative hypothesis (H_1)	assumed population parameter is <u>incorrect</u>	
	EXAMPLE: Conrobar	<ul style="list-style-type: none"> Days Absent – number of days employees are absent from work on sick/family leave Conrobar management claim that the average absenteeism rate is excessively high at over 2 days per employee. Set the null and alternative hypothesis NB: Remember that the analogy is: <ul style="list-style-type: none"> null = innocent alternative = guilty 	
	EXAMPLE: null hypothesis	<ul style="list-style-type: none"> Assume employees are <u>not</u> taking excessive leave until evidence demonstrates the contrary assume the avg employee absenteeism rate is <2 days Thus the null hypothesis can be written: $H_0: \mu \leq 2$ <p>NB: null hypothesis must always contain an equal sign</p>	


	<p>EXAMPLE: alternative hypothesis</p> <ul style="list-style-type: none"> statistical evidence is required to contradict the assumption contained in the null hypothesis Management's suspicion or contention is that employees overall have a mean that is > 2 days. \therefore the alternative hypothesis is: $H_0: \mu \leq 2$ 	<pre> graph TD A[Collect Random Sample n=48 employees, x̄=2.38 days] --> B{Is this consistent with Ho?} A --> C{Is this sufficiently different from Ho?} B --> D[Do NOT reject Ho] C --> E[Reject Ho in favour of H1] </pre> <p>Hypothesis test procedure enables us to choose which outcome</p>	
<p>Step 1: Setting up hypotheses</p>	<p>when are hypothesis tests used?</p>	<ul style="list-style-type: none"> in situations where we have: <ul style="list-style-type: none"> prior knowledge prior experience a standard a claim in these situations we have: <ul style="list-style-type: none"> some idea of value of population parameter being investigated, or some hypothesised value against which we can compare sample results 	
<p>Setting up (H_0) and (H_1) prior knowledge</p>	<p>EXAMPLE: assume prior knowledge is still current: set as H_0</p>	<p>Detailed survey in 2004 showed avg km travelled yearly per car was 14,500km - has there been any changes in 2015?</p> <ul style="list-style-type: none"> take random sample of cars want to know whether the mean is diff (has increase/decreased) <p>$H_0: \mu = 14,500$ no change: avg still 14,500km in 2015</p> <p>$H_1: \mu \neq 14,500$ change: avg not 14500km in 2015</p>	
<p>Setting up (H_0) and (H_1) prior experience</p>	<p>EXAMPLE: assume prior knowledge is still reliable: set as H_0</p>	<p>Past experience has shown no more than 3% of components from particular supplier are defective - new batch has arrived, test random sample of components from batch</p> <ul style="list-style-type: none"> could proportion (p > 3%) <p>$H_0: \mu \leq 3\%$ new batch has no more than 3% defective</p> <p>$H_1: \mu > 3\%$ new batch has more than 3% defective</p>	
<p>Setting up (</p>	<p>EXAMPLE:</p>	<p>Wish to test whet'</p>	

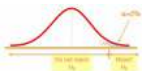

H_0) and (H_1) standard	normally set H_0 equal to std	<ul style="list-style-type: none"> take random sample of 48 we would write 2 hypotheses as: $H_0: \mu \geq 100\%$ employees on avg are meeting 100% std $H_1: \mu < 100\%$ employees on avg are failing to meet 100% std 	
Setting up (H_0) and (H_1) claim		Who is making claim and the seriousness of claim determines whether we: <ul style="list-style-type: none"> treat claim as null hypothesis (accept claim is true) until have evidence to contrary, or treat claim as alternative hypothesis (claim is false) 	
Setting up (H_0) and (H_1) claim	EXAMPLE:	Newspaper claims at least 40% of readers will see particular type of ad in its paper. Association of National Ads wants to check claim. <ul style="list-style-type: none"> adopt conservative approach + assume claim is true test claim with random sample of readers $H_0: \mu \geq 40\%$ at least 40% of readers ad $H_1: \mu < 40\%$ fewer than 40% see ad 	
	Common mistakes	<ul style="list-style-type: none"> equality part of hypotheses always appears in null hypothesis $H_0: p > 100$ $H_1: p \leq 100$ hypotheses are statements re population parameters <u>not</u> sample statistics $H_0: \bar{x} = 100$ $H_1: \bar{x} \neq 100$ 	
Step 2: Type of test	different types of tests	<ul style="list-style-type: none"> a hypothesis test re value of population mean (μ) or proportion (p) must take one of 3 forms: 	
Type of test	numerical data	$H_0: \mu \geq A$ $H_0: \mu \leq A$ $H_0: \mu = A$ $H_1: \mu < A$ $H_1: \mu > A$ $H_1: \mu \neq A$	
Type of test	categorical data	$H_0: p \geq A$ $H_0: p \leq A$ $H_0: p = A$	

		$H_1: p < A$	$H_1: p > A$	$H_1: p \neq A$	
--	--	--------------	--------------	-----------------	--

Type of test	one-tail test	<ul style="list-style-type: none"> sometimes, alternative hypothesis focuses on particular direction <div data-bbox="1010 134 1778 437"> <div> $H_0: \mu \geq 3$ $H_1: \mu < 3$ </div> <div> <p>⇒ This is a lower-tail test since the alternative hypothesis is focused on the lower tail below the mean of 3</p> </div> </div> <div> <div> $H_0: p \leq 8\%$ $H_1: p > 8\%$ </div> <div> <p>⇒ This is an upper-tail test since the alternative hypothesis is focused on the upper tail above the proportion of 8%</p> </div> </div>
--------------	---------------	---

		<p>∴ to set up H_0 and H_1 need to think about:</p> <ul style="list-style-type: none"> - Which error is the more serious (Type I) - The risk willing to run of making this error 	
	Level of significance (α)	<ul style="list-style-type: none"> • this is what hypothesis testing risk is referred to • specify maximum risk willing to accept in making Type I error • normally set at 5%, 10%, 2% or 1% • level of significance is max risk willing to accept in making Type I error (rejecting H_0 when shouldn't) 	
critical rejection region	critical rejection region	<ul style="list-style-type: none"> • consider if reject H_0 $H_0: p > 100$ $H_1: p \leq 100$ 	
critical rejection region		<ul style="list-style-type: none"> • Need to decide: <ul style="list-style-type: none"> - direction of test; and - specify(α) • Can next determine the critical (rejection) region = area in sampling distribution where we will reject H_0. • critical value(s) = border(s) of the critical region 	
	level of significance & critical (rejection) region		

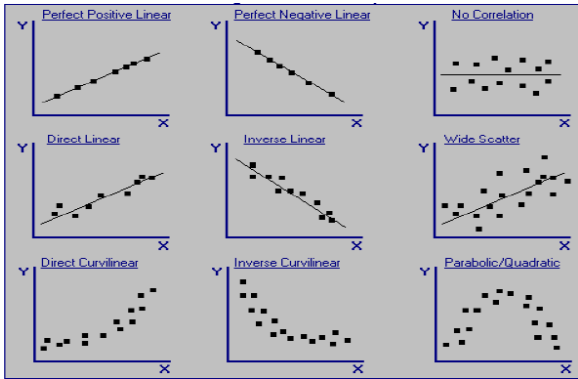
	critical value (CV)	<ul style="list-style-type: none"> can be determined in terms of Z or t-values EG. suppose $\alpha = 5\%$ (2-tail test) using Z 	
Step 4: Decision rule	Decision rule	<ul style="list-style-type: none"> if sample result falls beyond critical value = reject H_0 in favour of H_1 if calculate Z statistic from sample which is: $Z > +1.96$ or $Z < -1.96$ \therefore would reject H_0 	
Step 5: analyse sample	analyse sample	<ul style="list-style-type: none"> collect random sample calculate sample statistic \bar{x}, s for hypothesis tests re means \hat{p} for hypothesis tests re proportions calculate Z or t-statistic as appropriate 	
Step 6: conclusion	Conclusions	<ul style="list-style-type: none"> compare calculated test statistic (from step 5) to decision rule (from step 4) make decision: reject H_0 do not reject H_1 write conclusion re terms of problem solved 	
	EXAMPLE: Defectives	<p>Past experience show that no more than 3% of components from particular supplier are defective. New batch arrives - does this new batch of components have no more than 3% defective? If it has more, we will send it back to the supplier.</p> <ol style="list-style-type: none"> Take a random sample of components from batch to determine proportion of defective components. Then perform a hypothesis test to see if the population proportion in the batch could be greater than 3%. <p>Step 1: step up H_0 and H_1</p>	

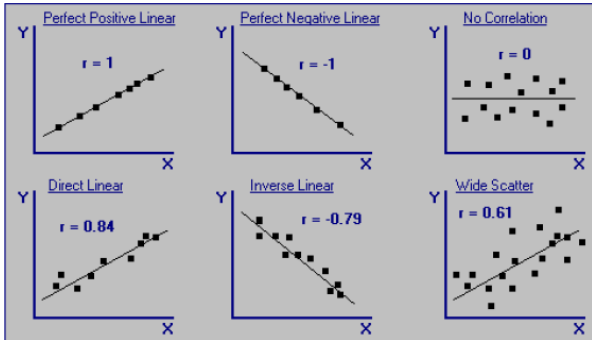
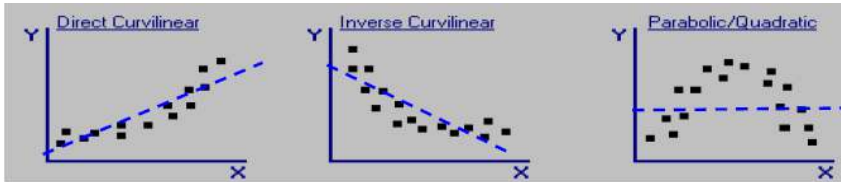
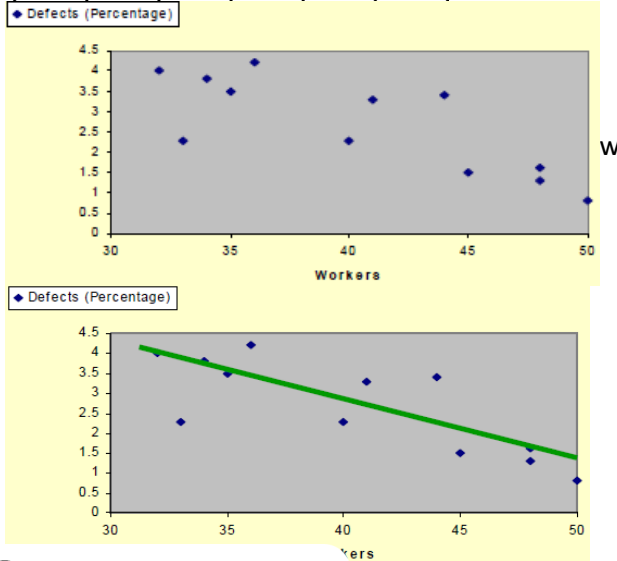
		<ul style="list-style-type: none"> • $H_0: p \leq 3\%$ = new batch has no more than 3% defective • $H_1: p > 3\%$ = new batch has more than 3% defective  <p>Step 2: decide on direction of test</p> <ul style="list-style-type: none"> • upper tail test <p>Step 3: Decision on α (lvl of significance)</p> <ul style="list-style-type: none"> • set α to be 5% • critical value (CV) of Z will be +1.645 <p>Step 4: decision rule (using critical values of Z or t)</p> <ul style="list-style-type: none"> • if proportion of defectives in sample result in Z static > 1.645 = reject H_0 <p>Step 5: sample (perform relevant calculations: Z or t-statistic)</p> <ul style="list-style-type: none"> • if test on random sample of 1000 components + 43 defective ...: <ul style="list-style-type: none"> - $n = 1000$ - $\hat{p} = \frac{43}{1000} = 0.043$ or 4.3% $\sigma \hat{p} = \sqrt{\frac{pq}{n}} = \sqrt{\frac{0.03(0.97)}{1000}} = 0.00539 \vee 0.539\%$ <ul style="list-style-type: none"> • find test statistic (z statistic) $z = \frac{\hat{p} - p}{\sigma \hat{p}} = \frac{0.043 - 0.03}{0.00539} = 2.41$ <ul style="list-style-type: none"> • test statistic = no. of std errors sample result from population parameter (assuming H_0 is true) • compare Z statistic to CV of Z <p>Step 6: conclusion (reject / not reject H_0 + answer question)</p> <ul style="list-style-type: none"> • Z statistic of 2.41 exceeds CV of Z (1.645) = reject H_0 (means sample proportion (evidence) fell in rejection region) • conclude at 5% level of significance - new batch has more than 3% defective • decision: send batch back 	
	<p>EXAMPLE: Tests on productivity at Conrobar</p>	<p>Are all staff on avg meeting 100% productivity std or better? (Use $\alpha = 5\%$)</p> <p>Step 1: step up H_0 and H_1</p> <ul style="list-style-type: none"> • $H_0: \mu \geq 100\%$ = staff meeting std • $H_1: \mu < 100\%$ = staff not meeting std  <p>Step 2: decide on</p>	

		<ul style="list-style-type: none"> lower tail test <p>Step 3: Decision on α (lvl of significance)</p> <ul style="list-style-type: none"> $\alpha = 5\%$ CV for t = -1.678 (from table) <p>Step 4: decision rule (using critical values of Z or t)</p> <ul style="list-style-type: none"> if sample results t static < 1.678 = reject H_0 <p>Step 5: sample (perform relevant calculations: Z or t-statistic)</p> <ul style="list-style-type: none"> $n = 48$ sample mean: $\bar{x} = 98.86$ sample SD: $s = 4.399$ Std error: $s \hat{x} = \frac{s}{\sqrt{n}} = \frac{4.399}{\sqrt{48}} = 0.635$ <ul style="list-style-type: none"> t-statistic: $t \hat{x} = \frac{\bar{x} - \mu}{s \hat{x}} = \frac{98.86 - 100}{0.635} = 1.795$ <p>Step 6: conclusion (reject / not reject H_0 + answer question)</p> <ul style="list-style-type: none"> t-statistic of -1.795 less than CV of t (-1.678) = reject H_0 sufficient evidence to conclude, at 5% lvl of significance, staff overall on avg not meeting productivity std of 100% 	
Different approach to hypothesis test	P-value	<p>the probability of getting test statistic more extreme than sample result, given null hypothesis (H_0) is true</p> <ul style="list-style-type: none"> derived from (Z or t) test statistic gives area in tail (or 2 tail if 2-tail test) beyond where sample result falls if p-value > α = do not reject H_0 if p-value < α = reject H_0 p-value <u>calculated by computer</u> 	
	EXAMPLES: Defective	<p>Step 5</p> <ul style="list-style-type: none"> p-value = 0.0080 or 0.8% <p>Step 6</p> <ul style="list-style-type: none"> p-value of 0.008 (or 0.8%) = < α of 0.05 (or 5%) \therefore reject H_0 Decision: the new batch of components has >3% defective 	
	EXAMPLES: Conrobar	<p>Step 5</p> <ul style="list-style-type: none"> p-value = 0.0395 or 3.95% 	

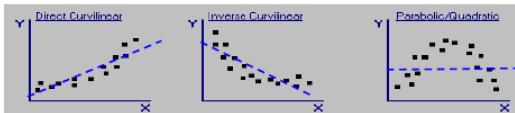
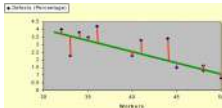
		Step 6 <ul style="list-style-type: none"> p-value of 3.95% = $< \alpha$ of 5% \therefore Reject v Decision: staff overall, on avg, are not meeting productivity std of 100%. 	
	Tests on productivity at Conrobar	Step 6: Conclusion at different α 's <ul style="list-style-type: none"> p-value = 0.0395 or 3.95%, tells us: <ul style="list-style-type: none"> would reject H_0 at an α of 10% (or $< 3.95\%$) would not reject H_0 at an α of 2% or 1% (or $> 3.95\%$) 	

Week 10 (Topic 10) - Chapter TEN (Simple Linear Regression)

	Regression & Correlation	<ul style="list-style-type: none"> relationships b/w variables enable us to explain nature of link b/w variables + how one variable affects the other build a model b/w variables for purpose of: <ul style="list-style-type: none"> - explanation - prediction - control - causality 	
Regression	Explanation	<ul style="list-style-type: none"> regression helps explain / understand variation in dependent variable do this by: finding other independent variables that relate to dependent variable wish to know: <ul style="list-style-type: none"> - direction of the relationship - strength of the relationship 	
	Prediction	<ul style="list-style-type: none"> make use of explanatory (independent) variable to predict likely outcome of dependent variable EG. knowing no. of customers fast food restaurant has may enable management to forecast sales 	
	Control	<ul style="list-style-type: none"> if we have some control over value of independent variable, this enables some form of control over dependent variable EG. varying advert expenditure up/down to certain extent = able to control movement in sales 	
	Causality	<ul style="list-style-type: none"> it is important to note that while regression models may establish association b/w variables - they do not necessarily establish causality 	
Concepts in regression & Correlation	Scatter diagram	<p>representation of possible relo b/w 2 variables</p> <ul style="list-style-type: none"> plots pairs of variables on scatter diagram to identify possible relationships verticle (y) axis) always contains dependent variable look for: <ul style="list-style-type: none"> - no relationships - linear relationships - non-linear relationships possible patterns 	

	<p>Correlation coefficient (r)</p>	<p>measure strength + direction or linear relationship b/w 2 variables</p> <ul style="list-style-type: none">• $-1 \leq r \leq 1$• r is relative measure of strength of linear relationship• closer r is to +1 or -1 = stronger linear relationship• r close to 0 = weaker / no linear relationship• sign of r provides direction of relationship 																										
	<p>non-linear relationship</p>	<ul style="list-style-type: none">• before interpreting r, scatter plot must be drawn• r associated with regression models would be poor indicators of actual strengths of each 																										
<p>EXAMPLE: % defects v. no. of workers</p>	<p>random sample of 12 shifts at manufacturer sampled and % of defective components recorded</p> <table border="1" data-bbox="611 873 1615 932"><tr><td>Defects</td><td>4</td><td>0.8</td><td>1.5</td><td>4.2</td><td>2.3</td><td>2.3</td><td>1.3</td><td>1.6</td><td>3.4</td><td>3.3</td><td>3.8</td><td>3.5</td></tr><tr><td>Workers</td><td>32</td><td>50</td><td>45</td><td>36</td><td>40</td><td>33</td><td>48</td><td>48</td><td>44</td><td>41</td><td>34</td><td>35</td></tr></table> <p>are defects & workers related?</p> <ul style="list-style-type: none">- negative linear relationship <p>correlation coefficient: $r = -0.788$</p> <ul style="list-style-type: none">- a moderate to strong negative linear relationship b/w defects & workers <p>Regression: line of best fit</p> <ul style="list-style-type: none">- estimated simple linear regression line: $\hat{y} = b_0 + b_1 x$ <p>\hat{y} = dependent variable x = independent variable b_0 = y-intercept (line cuts vertical axis) b_1 = slope of line</p>	Defects	4	0.8	1.5	4.2	2.3	2.3	1.3	1.6	3.4	3.3	3.8	3.5	Workers	32	50	45	36	40	33	48	48	44	41	34	35	
Defects	4	0.8	1.5	4.2	2.3	2.3	1.3	1.6	3.4	3.3	3.8	3.5																
Workers	32	50	45	36	40	33	48	48	44	41	34	35																

		<p>Regression Coefficients</p> <ul style="list-style-type: none"> - We explain using the defectives example: $\hat{y} = 8.480 - 0.144x$ - regression coefficients b_0 & b_1 can be interpreted in 3 ways: <ul style="list-style-type: none"> ▪ Geometrically (i.e. graphically) b_0 interpretation: on graph b_0 = where line cuts vertical axis EG. line cuts y axis at 8.48% b_{10} interpretation: on graph b_1 = slope of line EG. slop = -0.144 ▪ Algebraically (i.e. in equation form) b_0 interpretation: b_0 is value of y when x = 0 EG. y = 8.48 when x = 0 workers b_1 interpretation: b_1 is change in value of y when x changes by 1 EG. if x increases by 1 = y decreases by 0.144% ▪ Practically (i.e. practical interpretation) b_0 interpretation: b_0 not always useful interpretation as x = 0 may be outside range of x values used for regression equation EG. 8.48% of defects produced on avg when 0 workers = nonsensical b_1 interpretation: b_1 indicates impact on y from change in x EG. for each extra worker employed on a shift, on avg defectives decrease by 0.14% 	
Concepts in regression & Correlation	Regression	mathematical model of relationship	
	Simple linear regression	<ul style="list-style-type: none"> • involves 1 independent variable 	

	Multiple regression	<ul style="list-style-type: none"> involves 1+ independent variable to explain variation in dependent variable 	
	Non-linear relationship	<ul style="list-style-type: none"> before interpreting r, a scatter plot must be drawn r associated with regression models = poor indicators of actual strengths of each relo 	
how does model fit data?	Residuals	<ul style="list-style-type: none"> regression line does not perfectly fit data there will be variations (errors) b/w line + actual data points: $y - \hat{y}$ 	
how does model fit data?		<ul style="list-style-type: none"> need to obtain measures of these residuals & how well line fits with data measure variation around line = use <u>std error of estimate</u> (s_e) for how well line fits data = use <u>coefficient of determination</u> (r^2) 	
	Std error of estimates (s_e)	<ul style="list-style-type: none"> use to measure variation around line defectives v employees: $s_e = 0.758$ interpretation: estimate avg variation around regression line = 0.758% rough approx.. using empirical rule, could say max deviation from line be: $\pm (3 \times 0.758)$ or $\pm 2.27\%$ 	
	Coefficient of determination (r^2)	<ul style="list-style-type: none"> how well line fits data: $0 \leq r^2 \leq 1$ measures proportion of variation in 1 variable (y) dependent, explained by or attributable to variation in 2nd variable (x) independent provides absolute measure of strength of relo r^2 close to 1 = strong relo r^2 close to 0 = weak / non-existent relo normally expressed as % 	
	coefficient determination (r^2)	<ul style="list-style-type: none"> Defective v employees: $r^2 = 0.6209$ interpretation: <ul style="list-style-type: none"> approx.. 62% of variation in "defectives" explained by or attributed to variation in "no. of workers" in shift 	

		<ul style="list-style-type: none"> - remaining 38% variation = result of other factors not included in model (eg. time of shift, worker's experience etc) 	
	Use regression equation for estimation / prediction	<ul style="list-style-type: none"> • regression equation coefficient (b_0 and b_1) define nature of relationship b/w variations • regression equation also used for estimation / prediction <p>EG. if 50 workers in shift, estimated proportion of defectives would be:</p> $\hat{y} = b_0 + b_1 x$ <p>$\hat{y} \approx 8.48 - 0.144 * (50) \approx 1.28\%$</p> <p>This is a point estimate</p>	
	Prediction	when we use regression model with value of x contained in range of x values from sample	
	extrapolation	<p>when we use the model with value of x outside range</p> <ul style="list-style-type: none"> • should be used with caution as no guarantee same model holds outside original range of data 	
	EXAMPLE:	<p>Collected data over the last 10 years re annual expenditure on advertising & its total sales (all figures scaled for inflation).</p> <p>Develop a regression model and answer the following questions:</p> <ul style="list-style-type: none"> • How well does the model predict sales? • Interpret b_0 and b_1. • What would you estimate sales to be when \$1m is spent on advertising? 	