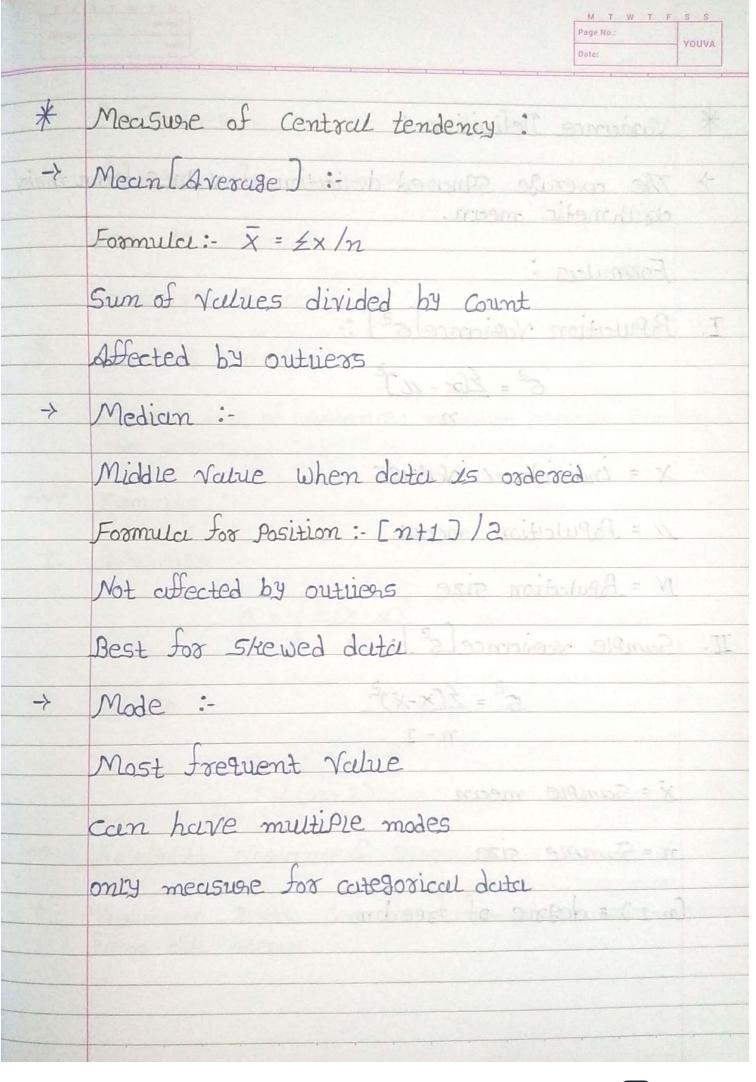
	Introduction To Statistics Page No. S Pa
*	What is Stats?
	Haubitihai Ita at Roll 1997
->	Stats is a science of Conecting, analyzing and
100	he again and place Property
*	What is Data?
- >	Facts or Pieces of information that can be measured
*	Types of Stats:
+	Descriptive stats:
	Organization of data.
->	Inferential Stats:
Die	Techniques whose used data
	that we have measured to form the condisions.
*	Sampling Techniques:
→	Simple Random Sampling:
	Every member of Population (N)
	have an equal chance of Jetting Selected in
beed!	Sample (n).
+0:	Stratified Sampling:
	Where we spirt the Population
	[N] into non-overlopping groups [strutu].
Name of Street or other Designation of the least of the l	

	Page Not: Page Not: YOUVA
	Systematic Sampling: Systematic Sampling: Systematic Sampling:
The Party of the P	a friedline fathous in some 2 a 31 Piole of
->	Convenience Sampling:- Surveying only those Searce who have knowledge of that Santicular domain.
*	Variables : 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
→	A variable is a Property that can take any Vul
*	
\rightarrow	Qualitative/Categorical Vasiables:
	Nominal: Categories with no order. [e.g., Colors, gender]
	Ordinal: - Catedories with natural order. [e.g., education level]
->.	Quantitative / Numerical Variables:
	an integral Printed to amond Louis and
	Discrete: Countable Values [e.g., number of chibren
neit	Continuous: - Any Values within a range leg, height,
	Copents I satisfy for 12012340 were atest [w]



		M T W T F	
		Date:	AOUVA
	Properties:	La Anna (P	4
	First any 2 Add at		
	ALWays non-negative		
		or SMI	
	Unit 15 squared		
	Sensitive to outliers	Engart M	5-
	stitute to outliers	77.0	4-
*	Standard Devication :-	2015	
	: 6th Resemblie	raikaM.	4-
->	5quare host of variance: measure and	age devis	tion
	From mean site of the state of	hoidy Teal	4
~~>			
	Formulas:	Projects/A	4
Ţ.	Population:	1	
Mark .	FAMILIA CONTRACTOR OF THE PARTY	JUIDIN 1	Took .
	$6 = \frac{1}{2}(x - u)^2$	= GOT	
	V N		
TI.	Sample:	Louves, Le	
-11-			
	$5 = \frac{\cancel{\xi}(x - x^{-})^{2}}{\cancel{(m-1)}}$	- Applies	
	$\sqrt{(n-1)}$	e i s	¥.
~~)	Key Point Vasiance & Standard	and O	2-
->			
	Variance gives you the average 59	mared dis	tunce

	M T W T F S S Page No.
	Date
+	Standard Deviction gives you a measure of
	squared in the same unit.
1	- Charles - Maritages - Marita
*	Five number summony:
->	Minimum . Comment of the last
	Minimum: 5 mallest value
->	Q1[First aucustive]: 25th Percentile
	-i mailtainen Acakras P - M
->	Medicin: 50th Percentile
moit	nively alternation appropriate to food accorded of
->	a3[Third aucostile]: 75th Percentile
- >	Maximum: Loogest Values
	Maximum . Lobdest Values
m	Calculations:
	IQR = Q3 - Q1 (N-X)
	Lower Lence = Q1 - 1.5[IQR]
	tower sence - at 1.31 tax
	Higher fence = Q3 + 1.5 [IQR]
4	(N-10) 2 = 7
*	Data Distribution:
4-	
/	Pattern of how data values age spred and their frequency of occurrence.
5.5	dequence of occushence.
	PENSAGE SAL MORE
	The state of the s

	M T W T F	
	Date:	YOUVA
my	Key Point / Component :-	1.
+	Center Joseph Office	4
+	Spread managed lastatests	4
+	shape	
	- State - Stat	
+	outiles 5	
Sive	Michaels declare man mem by the resultant	4
de	tions.	
*	Normal distribution:	
	- cumula-	
+	The normal distribution is a Probability distributhat occurs naturally in many Phenomena.	Hion
+ ~~	that occurs naturally in many Phenomena. Key characteristics:	tion
iaud	that occurs naturally in many Phenomena. Key characteristics:	tion
-	that occurs naturally in many Phenomena. Key characteristics: Beu-Shaped, Symmetric around the mean	ution.
hudi	that occurs naturally in many Phenomena. Key characteristics:	tion
À →	that occurs naturally in many Phenomena. Key characteristics: Bev- Shared, Symmetric around the mean 68% of data within 1 Standard deviation	Hion
→	that occurs naturally in many Phenomena. Key characteristics: Beu-Shaped, Symmetric around the mean 68% of data within 1 Standard deviation 95% of data within 2 Standard deviation	ution.
→ →	that occuss naturally in many Phenomena. Key characteristics: Ben-shared, Symmetric around the mean 68% of data within 1 Standard deviation 95% of data within 2 Standard deviation	Hion
→ → →	that occurs naturally in many Phenomena. Key characteristics: Ben-Shaped, Symmetric around the mean 68% of data within 1 standard deviation 95% of data within 2 standard deviation 95% of data within 3 standard deviation	tion
→ → →	that occurs naturally in many Phenomena. Key characteristics: Beu-shaped, Symmetric around the mean 68% of data within 1 standard deviation 95% of data within 2 standard deviation 99.7% of data within 3 standard deviation	Hion
+ + + +	that occurs naturally in many Phenomena. Key characteristics: Bev. Shared, Symmetric around the mean 68% of data within 1 Standard deviation 95% of data within 2 Standard deviation 99.7% of data within 3 Standard deviation Applications:	Hion
+ + + +	that occurs naturally in many Phenomena. Key characteristics: Beu-Shaped, Symmetric around the mean 68% of data within 1 Standard deviation 95% of data within 2 Standard deviation 99.7% of data within 3 Standard deviation Applications:	tion.
+ + + +	that occurs naturally in many Phenomena. Key characteristics: Bev. Shared, Symmetric around the mean 68% of data within 1 Standard deviation 95% of data within 2 Standard deviation 99.7% of data within 3 Standard deviation Applications:	tion.
+ + + +	that occurs naturally in many Phenomena. Key characteristics: Ben-Shaped, Symmetric around the mean 68% of data within 1 Standard deviation 95% of data within 2 Standard deviation 99.7% of data within 3 Standard deviation Applications: Nextural Phenomena [height, weight, Ia]	Hion
+ + + +	that occurs naturally in many Phenomena. Key characteristics: Ben-Shaped, Symmetric around the mean 68% of data within 1 Standard deviation 95% of data within 2 Standard deviation 99.7% of data within 3 Standard deviation Applications: Nextural Phenomena [height, weight, Ia]	tion

Page No.: Date:	ADDAT
Financial markets	1
Queuity control	-
Statistical inference	1
Z-5000e:-	-
Measures distance from moun in student	4
Donulle.	14.
Z = X - II	4
6	
Used for compasing values from different distrib	utlor
The state of the property and the second	
Finance:	-
Standard Normal Distribution:	-
restantion hardwarts e middle whole is a in a re-	1
	and the same
Used as selence distribution	
Enables Probability calculations using z-table	
	Queuity control Statistical inference. Z-5core: Measures distance from mean in standard delions. Formula: Z = X-11 G Used for comparing Values from different distrib Common in statistics, Education [test score] Finance. Standard Normal Distribution: Mean = 0, Standard deviation = 1

	18 T W T F	8 8
	Date;	AOUAY
		-
¥	Normalization:	7
1	1 OD MICHELL ZCITIONL.	
1	Mallala o Mi	12
->	Methods: - Min - Max, z-score, decimal scaling	
1	min Man formulas a visitadis i	1-
->	Min-Mase Formula: - X - min	
	Xmax-Xmin	10
->	Used in machine learning, duta Pregocessing	
,	O to A ' live on Giorgia Gradation de Corio	*
->	Puts variables on same scale	
	Countrible Malzins	1-
*		
1	Ske wed Distributions:	1-
~~	Positive skew:	1-
_\	Mean > meadian	
	Mean / meadlest	
->	Long tale sight	- 6
->	Example: income distribution	4
The second		
m	Negative skew:	K-10
	D=0	
->	Mean & meadian	
	Medan - coth	
->	long tail loft	
->	Example: - exam scores	
	The state of the s	

		M T W 1	F 3 3
		Date:	YOUYA
*	Continuous Distribution:	s it were to	
->	Infinite Possible Values	: Platfor	1 4
->	Example: height, time	nin-Max	
->	Uses integration for Probability.	r al les	
*	Discrete Distribution:		
->	Countable values		
->	Example: - dice volls, counts		*
→	Uses Summertion for Probability	Switing	+
*	Vnisozm Distribution:	of Errol	4
->	Every outcome equally likely		
->	Formula: $f(x) = 1$ $b-a$	Widowall.	-fron
->	Mean: - atb	int pent	
->	Vasiance: - [b-a] ²		

	Page No.: YOUVA
->	Example: dice 70115, Irandom number generators.
*	Exponentical Distribution:
->	Models time between events
+	Formula: $f(x) = \lambda e^{-\lambda n}$ for $x \ge 0$
->	Mean: 1/2
->	Variance: 1/2 (9-1) 911-1911
->	Memoryless Property
- >	Example: equipment failure times, customes assivat
*	A Rose exemts in tend intend
->	
->	Bernouui Distribution:
	Bernouui Distribution:
->	Bernouwi Distribution: Single trical, Success / fairure $P(x=1) = P$, $P(x=0) = 1-P$ Mean: P
-> ->	Bernouwi Distribution: Single trial, Success / fallowe $P(x=1) = P$, $P(x=0) = 1-P$ Mean: P Yariance: $P(1-P)$
	Bernouwi Distribution: Single tricul, Success / fairure $P(x=1) = P$, $P(x=0) = 1-P$ Mean: P

	M T W T	FSS
	Date	YOUVA
d.		
*	Binomical Distribution:	1
->		
	n independent Bernouui trials	
->	Formula:	1
		1
	$J^{2}(X=k) = m p^{2} (1-p)^{-1/2}$	
	OZX TOT IN E CO) + -: Wilsonsol	14
->	Mean: - Mp np	
		- da
->	Variance: - np (1-p)	6
->		
	Example: - number of heads in n coin files	ďω
N/z	Promote: equipment fulnice times alamine	
*	Poisson Distribution:	
->	0 00 00/0 1 : (:)	
	Ruse events in fixed interval	16
->	Formula:	A
		4-
	$P(X=K) = 2^{K}e^{-2}$	
	9-K-1(0=x)9 9=(1=x)9	(-)
	Menn: P	
->	Mean = Variance = 2	
	Variance: P(2-P)	4
->	Example: - Website visits per hour	1
	Frankle : Single (am till	

Probability: P(E) = favorable outcomes / total outcomes P(Complement) = 1-P(E) Hypothesis Testing : Particular hypothesis test is a method of statistical inference used to decide whether the deta at had is sufficient to support a proticular hypothesis. Hypothesis testing alows as to make Probabilisticatement about population recurrences; Null Hypothesis [Ho]:- The null hypothesis assums that there is no significant recurrenship or effect blue two regions to static study or the assumption of no effect up proven otherwise. Herenate Hypothesis[Ho or H1]:- H is a statement, that contradicts the NH & claims these is significance effect or recurrens.		M T W T F	YOUVA
P(E) = favorable outcomes / total outcomes > O = P(E) = 1 > P(complement) = 1-P(E) * HyPothesis Testind: > A statical Aylothesis test is a method of statistical inference used to decide whether the data at had is sufficient to support a Proxiticular hypothesis. > HyPothesis testing allows as to make Probabilist statement about population Recuneres. * Null HyPothesis [Ho]: > The null hyPothesis assums that there is no significant rectionship or effect blo two decides > It solves as a statement point for HT & represent to state the sumption of no effect up proven otherwise. * Alternate HyPothesis [Ho or H1]:- > It is a statement, that contradicts the NH &		Date)	
A Statical hypothesis test is a method of statistical infermence used to decide whether the deta at had is sufficient to support a proticular hypothesis. Hypothesis testing allows as to make probabilisticatement about population procumeters. Nun Hypothesis [Ho]: The nun hypothesis assums that there is no significant pertionship or effect blue two herical to static que' or the assumption of no effect up proven otherwise. Alternate Hypothesis [Ho or H1]:- H is a statement, that contradicts the NH &	*	Probability:	
* P(complement) = 1-P(E) * HyPothesis Testing: -> A statical hyPothesis test is a method of statistical inference used to decide whether the data at had is sufficient to support a Proficular hypothesis. -> HyPothesis testing allows as to make Probabilist gatement about population recurrences. * Null HyPothesis [Ho]: -> The null hyPothesis assums that there is no significant relationship or effect blo two related to static sun or the assumption of no effect up proven otherwise. * Alternate HyPothesis[Ho or H1]:- -> It is a statement, that contradicts the NH &	->	P(E) = favorable outcomes / total outcome	5 X
Hypothesis Testing: A Statical hypothesis test is a method of statistical inference used to decide whether the data at had is sufficient to support a particular hypothesis. Hypothesis testing allows as to make Probabilist statement about population recurrences. Null Hypothesis [Ho]: The null hypothesis assums that there is no significant relationship or effect blue two regions. The serves as a statement point for HT & represented in State star statement of no effect we prove otherwise. Alternate Hypothesis[Ho or H1]:- H is a statement, that contradicts the NH &	+	0 \(\text{P(E)} \(\text{1} \)	.64
A Statical hypothesis test is a method of Statistical informence used to decide whether the data at had is sufficient to support a Prosticular hypothesis. † Hypothesis testing alows as to make Probabilist statement about population programeders. Null Hypothesis [Ho]:- The null hypothesis assums that there is no significant recrionship or effect blook two weight The saves as a stanting point for HT & gerse to Static row or the assumption of no effect we proven otherwise. # Alternate Hypothesis [Ho or H1]:- The is a statement, that contradicts the NH &	4	P (complement) = 1-P(E)	-8
Statistical inference used to decide whether the data at had is sufficient to support a propositional hypothesis. Hypothesis testing alows as to make probabilistic statement about population procunetess. * Null Hypothesis [Ho]:- The null hypothesis assums that there is no significant gentionship or effect blue two statements is static two or the assumption of no effect we prove otherwise. * Alternate Hypothesis [Ho or H1]:- The is a statement, that contradicts the NH &	*	HyPothesis Testing:	.8
# Null Hypothesis [Ho]:- The null hypothesis assums that there is no significant neutionship or effect blu two kinds The sorves as a starting point for HT & nepres to static run or the assumption of no effect us proven otherwise. # Alternate Hypothesis [Ho or H1]:- The is a statement, that contradicts the NH &	->	Statistical inference used to decide whether duta at had is sufficient to support a	the
The new hypothesis assums that there is no significant reactionship or effect blw two region of It sorves as a starting point for HT of general to Static grow or the assumption of no effect we proven otherwise. * Alternate Hypothesis[Ho or H1]:-	7		
JE SERVES as a Statement, that contradicts the NH &	*	Nun Hypothesis [Ho]:	9
* Static que or the assumption of no effect we proven otherwise. * Alternate Hypothesis[Ho or H1]:- The is a statement, that contradicts the NH &	>		
* Alternate Hypothesis[Ho or H1]:-	->	Proven otherwise.	t until
It is a statement, that contradicts the NH &	*	Alternate Hypothesis[Ho or H1]:-	
	->	It is a statement, that contradicts the NH	æ



		M T W T P Page No.: Date:	Youva
	ship.	Andoss	
*	Rejection Region Method:	(3)9	4.
1.	Ho & Ha	0220	6
2.	& -> Value and-1-14	9m N 9	
3.	cussumptions English Pipa	HAPPH	-
4.	decide test -> z-test, r-test	to 1	-(-
5.	value	1-10	
6.	Test Conduct		
7.	Reject / Accept	Carless	*
8.	State hesuits	11 Vivil	7
*	T-test :- tolk a dilamina +	a of	+
->-	Criticology	between	
->	3 TYPES :-	Person	
1-	200 1 00 000		
- 17	The state of the s	ai g	4

	M T W T F S S
	Date: YOUVA
>	Companes the mean of a single sample to a
	known
2.	Independent The sample t-test:
1	The samples from two independent groups are
-	companed to determine if the means of the associa
	ted Populations are significantly different.
	E = Now total X Column total
	$t = \bar{x}_1 - \bar{x}_2$
	512 + 522
	$\sqrt{n_2}$ n_3
	(2-1) x (c-2)
3.	Palsed t-Test:
1	compains the means of two hercited groups.
7	Compains the metris we can start and the
	# = <u>d</u>
	5d/Vn
	# Xustq515 :
	d = 2 meurs différence
	to be transfer to the transfer of the training of the
	5d = 5td difference.

* chi - square test :-

It is a Statical Procedure for determine the -> difference between observed and experted duta.

		M T W 1 Page No.:	100
		Delec	YOUNG
	2 = Z (0-E)2		
	E	TY locar	
	0 = observed frequency		
ati-	E = Expected frequency		
	F = Row total x column to	0.00	
	around total		
	$f = (9-1) \times (c-1)$		
	9 = no of 90ws		
	c = no of columns	aum c	
*	Kuntosis:		
→	Kurtosis is a measure of the tailedne distribution.		
>	3 Types of Kuntosis	2-14	16
	THE ENGLANCE AND DEVENTED AND ADDRESS OF THE PARTY OF THE	70 /T	

		M T W T F	s s Youva
		Date:	10072
1.	Mesokustic:		
	-> Tails use similar to ND		
	> distribution with kustosis = 3		
	Eg. Standard Normal distribution		
2.	Leptokustic :-		
	-> Heavy tails [with more / extreme or	utuens]	
	-> distribution with kustosis > 3		
	Eg. t-distribution with very small of.		
3,	Platy Kuntic:		
	> 1ight toins [fewer extrane outliers].		
	-> distribution with Kustosis <3		
	Eg. Uniform dest		,
			· · · · · · · · · · · · · · · · · · ·