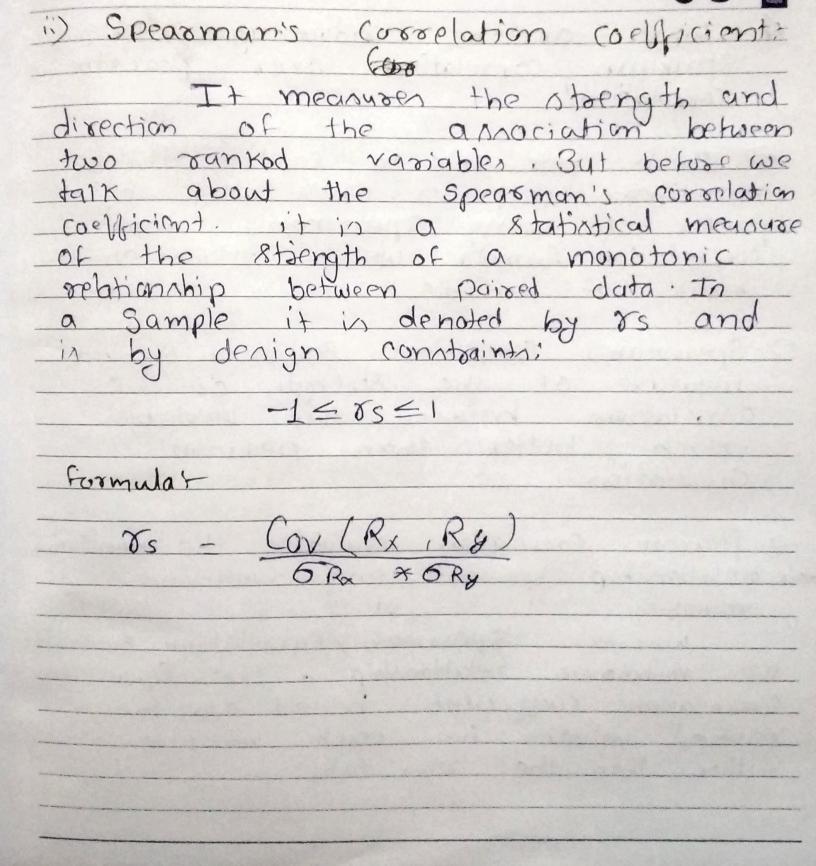
1. What in the definition of
covasionce? Create formula for it.
- Covariance in the measure of the
relationship between two dandom
Manaples and to what and
they change together.
- ather 1000 1000
can say, it defines changes between two variables, such that change in one variable is equal to change
two variables, such that change in
one variable is equal to change
in another Variables, the
There are two Kinds of Covariance:
1) Positive Covariance
ii) negative Covariance.
P
Formula !-
For Population:-
Cov (n,y) = = (x-x)(y-y)
Cov(n,y) = Z(x-x)(y-y)
6 - 6 1
for Sample:
Cov (n,y)
-0 · (· · / 3 /

2. What makes correlations better than Covariance? -> Covariance and correlation are two terms which are enactly opposite to each other, they both are used in statistics and regression analysis. Covariance shows on how the two variance variable vary from each other wherean correlation shows un the relationship between the two variables and how are they related. Correlation and covariance are very closely related to each other and yet they differ a lot covariance delines the type of infraction but correlation defines not the type only but also the strength of this reason of sunday correlation is often termed as the sperial cone of Covariance. However if one must choose between the two most will choose correlation an it remain unappected by the

Changes in dimensions, location and scale. Also, since it range between -1 to +1 it is useful to draw companisions between variables across domains.
3. Explain the process on well as Pearson and spearman correlation.
Correlation is a degree to which two variables are linearly related. This is an important step in bi-variant data analysis. In the broadent sense correlation is actually any statistical relationship whether casual or not.
a statistical measure of the strength of the relationship between the relative movement of two variable. The value range between -1.0 and 1.0 A correlation of -1.0 shows negative correlation, while a correlation of 1.0 Shows a positive correlation. A

consolation of 0.0 shows no movement of two variable There are two important Correlation Coefficient :-1. Pearson Correlation Coefficient in a measure of the strength of linear amoriation between two a pearson coefficient draw a line of bent fit through the data of two variable, and the pearton correlation coefficient, or indicates how for away all data points are to this line of best fit. Formilla -5 (niv) = Cov (niv) 6n ×6y



4. What are the advantage over
4. What are the advantage over Sprarman correlation over pearson
correlation?
-> It is not able to capture the
prefer of using Spearman Rank
prefer of using Spearman Rank
Tank of variable.
-> Speasman's Coefficient Res as a significant measure of the strength of the
amociation between two variables
much better than peamon
Correlation
-> Pearnon Correlation evaluates the mondanic mas relationship between two Continuous rapiables.
whereas Spearman correlation and to
the monotonic relationship. The Spearman
ranked values has each variable
rather than the raw data
THE MAIL

5. Describe	the	centra	d Limi	+
theorem.				
7				
shape of the distribution of always be	Samp	regardless ation distributed means imately	m will	
-> 9f we take Population dintri those sample dintribution.	bution wou	mple of the r	neam of	E mal
become more size increase	norma	Sample I an i	m eans	will
Sample distrib approximatly 8:Ze is	normal	it The	11 be Saw 30-	nple