= Mycompanion 43304

Aire + Principal Component Analysis

Finding principal components, variance and S.D. Colorlations of principal components (using R)

* Theory
Principal Component Analysis:

It is a rethood of extracting important variables (in four of components) from a large set of variables available in a dataset. It extracts low directional set of beatures from a high directional dataset with a nature to capture as much information as possible. It is useful for dealing with 3 on higher directional data.

It is always performed on a symmetric correlation on covariance mothers. This means that the nation should be nuronic and have standardized data.

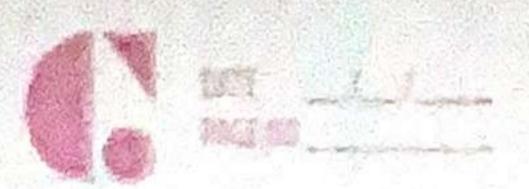
The first principal component gives The direction of the manimum spread of data. The second gives the direction of the manimum spread for to the first direction. Each solvers gives a direction.

Normalieation:

The principal components are supplied with portralized version of original predictors. This is because the original predictors may have different scales.

large loading for variables with high variance. In turn this will lead to insanely will lead to dependence of a principal component on the variable with high variance, which is undesirable.

Variance and Co-variance are a reasure of spread spread of a set of points around their center of mass (man).



ANY COMPANION
Variance is resource of the deviation from the reas for
paints in one dimension.
: Variance > 1 Standard Deviation
Lo-variance is measure of how much early of the dimensions
vary from the moan wirit each other.
: Covariance (X, V) = £ (Xi - X) (Ti - T)
(n-17
Standard Deviation:
Heavure of spread of the data points.
i o o I (ni - near)
Complement The thin assistance to the Hill and
conclusion - Thus in this assignment, we studied and implemented the concept of PEA on Big Mont dataset using
R. Marie and the state of the s