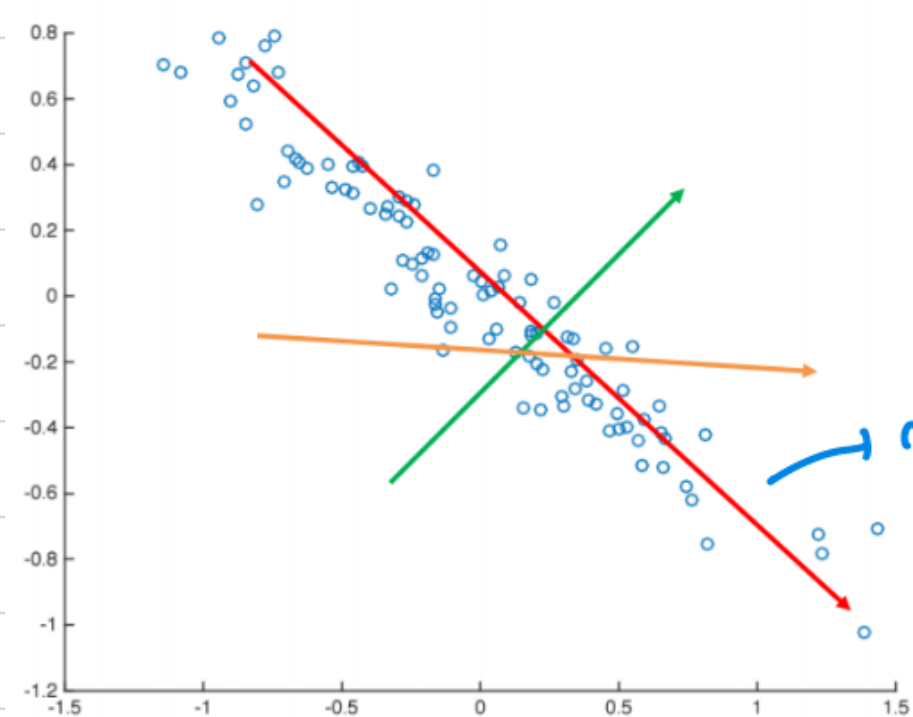# Principal Component Analysis

↳ Method for finding direction in high dimensional data that contain information } Dimension Reduction

↳ · unsupervised
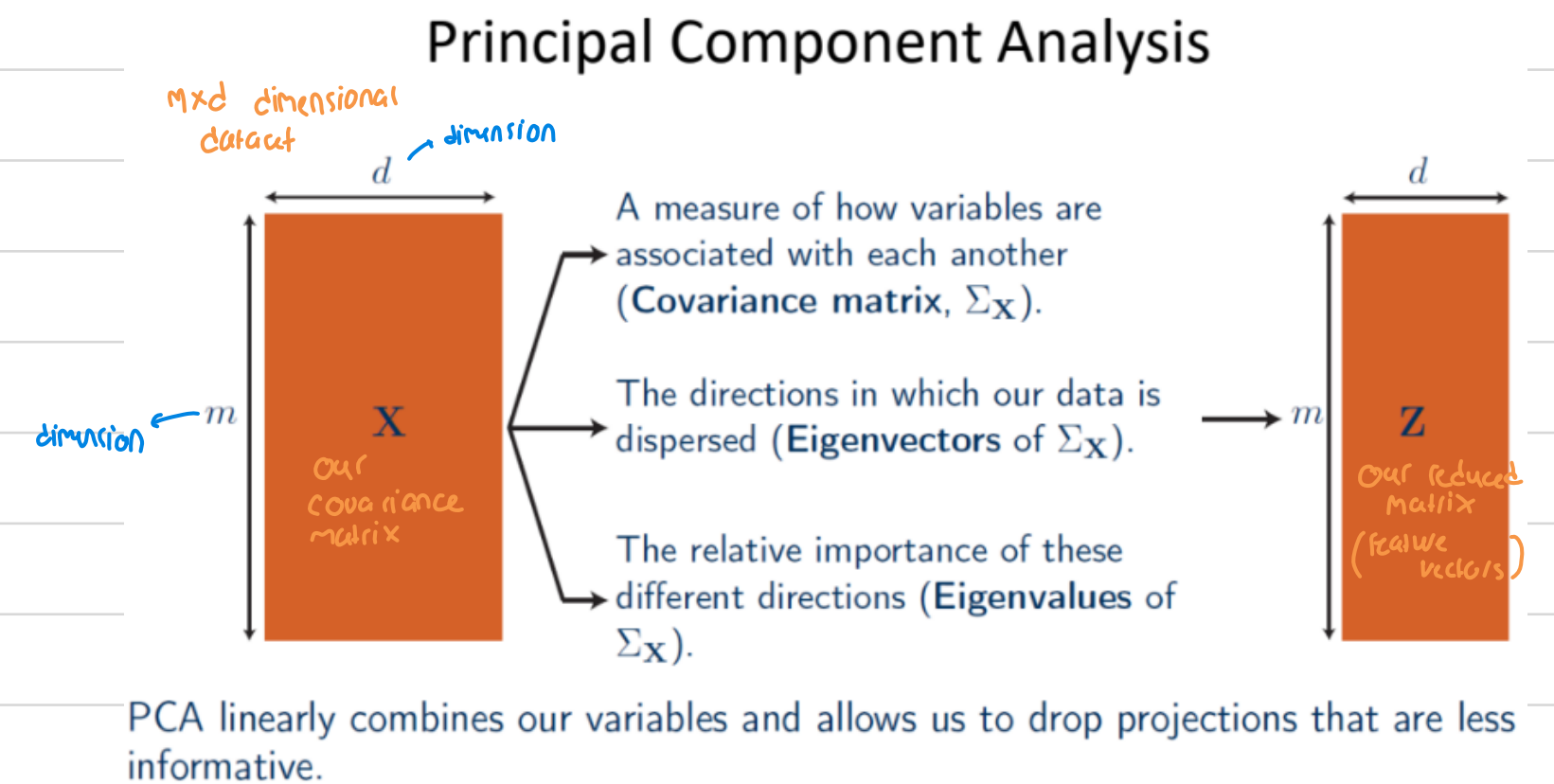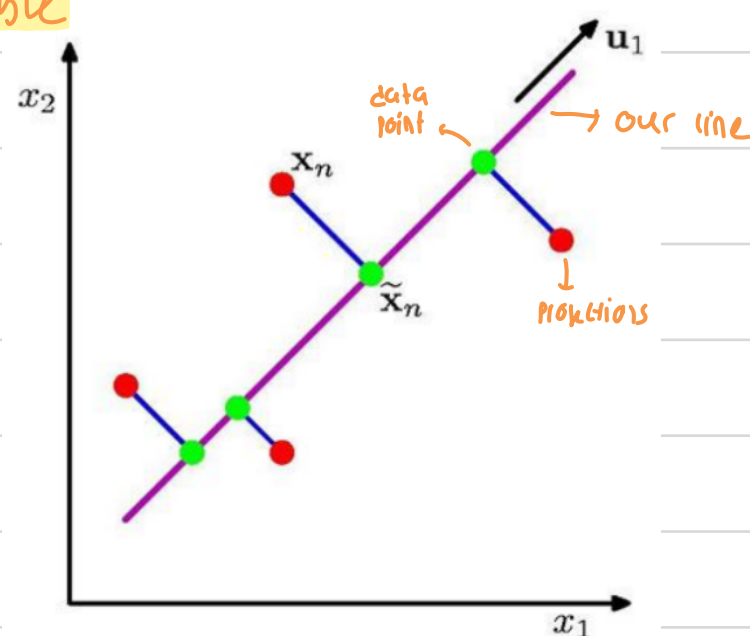· non-parametric

↳ examine interrelations among set of variables

↳ Goal :- ○ reduce the dimension of dataset while preserving the most important pattern or relationships between variables

○ by finding new set of variables containing most sample's information



## Principal Component Analysis

Mxd dimensional dataset

dimension

$d$

dimension ← $m$

**X**
our covariance matrix

A measure of how variables are associated with each another (**Covariance matrix**, $\Sigma_X$).

The directions in which our data is dispersed (**Eigenvectors** of $\Sigma_X$).

The relative importance of these different directions (**Eigenvalues** of $\Sigma_X$).

$d$

$m$

**Z**
Our reduced matrix (feature vectors)

PCA linearly combines our variables and allows us to drop projections that are less informative.

Which direction gives you the <u>maximum variance?</u>
Green / Orange / Red?

max variance

Choose a line that fits the data, so that :-

↳ the projection maximizes variance/variable of projected data (purple line)

↳ minimizes mean squared distances between :-
· data point &
· Projections (blue line)



① Eigenvectors & Eigenvalues :-
    ↳ directions of our spread w/ most info.
    ↳ each variable/feature has an eigenvectors &
    ↳ each eigenvectors have Eigenvalues :-
        ↳ How much information carried
    ↳ Rank our eigenvectors by eigenvalues, from low to high
        $PC_1$ = highest, $PC_2$ = 2nd highest, . . .

② Feature Vectors :-
    ↳ Matrix that has eigenvectors of components we want to keep
    ↳ where dimension reduction happen
    ↳ keep or discard eigenvector w/ lower significance

③ Recast the Data along Principal Component Axes
    ↳ to recast our data from the original axes to Principal Component using our feature vector
    => Final Data Set = Feature Vector$^T$ * Original Data Set$^T$
                         transposed          transposed
                         Feature vector      og Data set