## PCA and VAE

Sunday, 8 February 2026   09:46

# Dimensionality Reduction techs.

## PCA

- a statistical technique used to reduce the dimensionality of large data sets while preserving as much variance as possible.
  - reduces number of features in a data set while keeping the most important information
  - it changes complex data sets by transforming correlated features into smaller sets of uncorrelated components.

- it helps us remove redundancy and improve computational efficiency while making the data easier to visualise.

- It uses linear algebra to transform data into principal components
- it does this by calculating eigenvectors (directions) and eigenvalues (importance) from the covariance matrix

- Step 1  Standardise the data
- Step 2  Calculate Covariance matrix
- Step 3  find the principal components.
- Step 4  Pick the top Directions and Transform Data
- Can be done in python using Sk learn

### + ve

1. multicollinearity Handling : Creates NEW uncorrelated variables to address issues when original features are highly correlated.

2. Noise Reduction : reduces components with low variance plus increases data clarity.

3. Data Compression : Reduces data size

4. Outlier detection : identifies outliers.

### - ve

1. Interpretation Challenges : Principal components are combinations to can be hard to explain

2. Data Scaling sensitivity: Requires proper scaling of data or results will be misleading.

3. Information Loss : may lead to loss if too few components are kept.