

## 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 visualize.
  - It uses linear algebra to transform data into principal components
  - it does this by calculating eigen vectors (directions) and eigen values (importance) from the covariance matrix.
- Step 1 Standardize 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 sklearn

+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 thus increasing data clarity.
3. Data compression: Reduces data size
4. Outlier detection: identifies outliers.

-ve

1. Interpretation challenges: Principal components are combinations so 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.
4. Assumption of linearity: May struggle with non-linear data
5. Computational complexity: Can be slow and resource intensive on large data sets
6. Risk of overfitting.

Variational auto encoders.

are generative models that learn a smooth, probabilistic latent space.

- VAE's capture the underlying structure of a dataset and produce outputs that closely resemble the original data.

- learns a continuous latent representation
- Enable controlled and meaningful data generation
- widely used in image synthesis, anomaly detection and representation learning

• Step 1 Encoder (understand the input)

• Step 2 Latent space (adding some randomness)

• Step 3 Decoder (reconstructing / creating new dataset)

+ve

- Generative modelling

- Anomaly detection

- Data imputation and denoising