

# Principal Component Analysis (PCA)

---

PCA is a linear dimensionality reduction technique that uses Singular Value Decomposition (SVD) to project high-dimensional data onto a lower-dimensional space while preserving as much variance as possible.

## Key Concepts

---

- **Dimensionality Reduction:** Reduces the number of features in a dataset, simplifying analysis and improving model performance.
- **Variance Maximization:** The first principal component (PC) captures the direction of maximum variance in the data. Subsequent PCs are orthogonal to the preceding ones and capture progressively less variance.
- **Singular Value Decomposition (SVD):** A matrix factorization technique used by PCA to find the principal components.
- **Explained Variance Ratio:** The proportion of total variance captured by each principal component. This helps determine the number of PCs to retain.
- **Orthogonality:** Principal components are mutually orthogonal (uncorrelated), ensuring independence of the new features.

## PCA Algorithm

---

PCA involves these steps:

1. **Center the data:** Subtract the mean of each feature from the corresponding feature values.
2. **Compute the covariance matrix:** This matrix describes the relationships between the features.
3. **Perform SVD on the covariance matrix:** This decomposes the matrix into three matrices: (U), (S), and ( $V^T$ ). The columns of ( $V^T$ ) represent the principal components.
4. **Select the top k principal components:** Choose the number of PCs based on the desired variance explained or a predetermined number.
5. **Project the data:** Transform the original data onto the selected principal components to obtain the reduced-dimensional representation.

## Mathematical Formulation

---

The projection of the data matrix (X) onto the first (k) principal components ((W)) is given by:

$$[X_{\text{proj}}] = XW$$

where (W) is a matrix whose columns are the top (k) eigenvectors of the covariance matrix.

# Implementation with Scikit-Learn

---

```
from sklearn.decomposition import PCA
import numpy as np

# Sample data (replace with your data)
X = np.random.rand(100, 10)

# Apply PCA to reduce to 2 dimensions
pca = PCA(n_components=2)
X_reduced = pca.fit_transform(X)

# Access principal components
principal_components = pca.components_

# Explained variance ratio
explained_variance_ratio = pca.explained_variance_ratio_
```

The `explained_variance_ratio_` attribute shows the proportion of variance explained by each PC. Cumulative sum can be used to determine the number of PCs needed to reach a certain variance threshold (e.g., 95%).

## Choosing the Number of Dimensions

---

The number of principal components to retain can be determined in several ways:

- **Variance Explained:** Select the number of PCs that capture a sufficient percentage (e.g., 95%) of the total variance.
- **Scree Plot:** A plot of explained variance vs. number of PCs; the "elbow" point often indicates a suitable number of PCs.
- **Predetermined Number:** Choose a fixed number of PCs based on computational constraints or desired level of dimensionality reduction.

## Advanced Techniques

---

- **Incremental PCA:** Handles large datasets that don't fit in memory by processing them in batches.
- **Randomized PCA:** A faster approximation of PCA, especially useful for high-dimensional data.
- **Kernel PCA:** Extends PCA to non-linear relationships using kernel methods.

## Conclusion

---

Principal Component Analysis is a powerful dimensionality reduction technique widely used in machine learning for feature extraction, noise reduction, data visualization, and improving the efficiency of various machine learning algorithms. Its ability to preserve variance while reducing dimensionality makes it invaluable for handling high-dimensional datasets.