

## Lesson 31 Summary: PCA

Data Science with Python – Key Concepts

Data Science Program

## Principal Component Analysis

### Core Idea

Find directions of  
max variance  
Orthogonal components

### Variance Explained

Cumulative ratio  
Choose n for 95%  
Scree plot / elbow

### Use Cases

Dimensionality reduction  
Visualization  
Noise reduction

### Prerequisites

Scale features first (StandardScaler)  
Sensitive to outliers

### Finance Applications

Factor extraction | Risk decomposition  
Portfolio analysis

`PCA(n_components=2).fit_transform(X_scaled) | pca.explained_variance_ratio_`

---

PCA finds the directions of maximum variance

## Dimensionality reduction:

- **Find:** Directions with most variance
- **Project:** Data onto these directions
- **Result:** Fewer dimensions, preserved structure

---

Principal components are orthogonal to each other

## Choosing number of components:

- **Cumulative ratio:** Sum of explained variance
- **Rule of thumb:** Keep 95% of variance
- **Scree plot:** Look for elbow

---

First few PCs often capture most of the variance

## Why standardize:

- **Problem:** PCA sensitive to scale
- **Solution:** StandardScaler before PCA
- **Result:** Equal contribution from all features

---

Always scale features before PCA

## Understanding PCs:

- **Loadings:** Weights of original features
- **High loadings:** Important for that PC
- **Sign:** Direction of relationship

---

Access loadings via `pca.components_`

## Common uses:

- **2D scatter:** PC1 vs PC2
- **Color by label:** See cluster separation
- **Biplot:** Show feature contributions

---

PCA is a powerful visualization tool

## When PCA struggles:

- **Non-linear:** Only captures linear relationships
- **Outliers:** Can distort principal components
- **Interpretation:** PCs may not be meaningful

---

Consider t-SNE or UMAP for non-linear data

## Use cases:

- **Factor extraction:** Find market factors
- **Risk decomposition:** Identify risk sources
- **Portfolio compression:** Reduce dimensionality

---

PCA is widely used in quantitative finance

### Essential Commands:

Task	Code
Scale data	<code>StandardScaler().fit_transform(X)</code>
Fit PCA	<code>PCA(n_components=2).fit(X_scaled)</code>
Transform	<code>pca.transform(X_scaled)</code>
Variance ratio	<code>pca.explained_variance_ratio_</code>
Loadings	<code>pca.components_</code>

---

**PCA is a fundamental unsupervised technique**