

Here are the completed answers for all the questions on Principal Component Analysis (PCA).

# Pca Interview Questions - Theory Questions

## Question 1

### What is Principal Component Analysis (PCA)?

#### Theory

Principal Component Analysis (PCA) is a fundamental, unsupervised linear dimensionality reduction technique. Its primary goal is to transform a dataset with a large number of correlated variables into a smaller set of new, uncorrelated variables, called **principal components**. These components are ordered by the amount of original variance they capture, with the first component capturing the most, the second capturing the second most, and so on. By retaining only the first few principal components, PCA reduces the dimensionality of the data while preserving as much of its essential information (variance) as possible.

#### Explanation

PCA achieves this transformation by finding a new set of orthogonal axes (the principal components) that are aligned with the directions of maximum variance in the data.

1. **Standardize the Data:** PCA is sensitive to the scale of the features. The first step is to standardize the data to have a mean of 0 and a standard deviation of 1.
2. **Compute Covariance Matrix:** It calculates the covariance matrix to understand how the variables in the original dataset relate to each other.
3. **Find Eigenvectors and Eigenvalues:** It performs an eigen-decomposition of the covariance matrix to find the eigenvectors (which represent the directions of the principal components) and the eigenvalues (which represent the magnitude of variance captured by each component).
4. **Select Principal Components:** The components are sorted by their corresponding eigenvalues in descending order. The top k components that capture a desired amount of variance (e.g., 95%) are selected.
5. **Project Data:** The original data is projected onto the new subspace defined by the selected principal components.

#### Use Cases

- **Dimensionality Reduction:** Reducing the number of features in a dataset to combat the curse of dimensionality, speed up model training, and reduce storage requirements.
  - **Data Visualization:** Projecting high-dimensional data onto 2 or 3 principal components to visualize its structure in a 2D or 3D scatter plot.
  - **Noise Reduction:** By discarding components with low variance (which often correspond to noise), PCA can be used to clean up a dataset.
  - **Feature Engineering:** The principal components themselves can be used as new, decorrelated features for supervised learning models.
- 

## Question 2

Can you explain the concept of eigenvalues and eigenvectors in PCA?

### Theory

Eigenvectors and eigenvalues are the mathematical core of PCA. They are derived from the covariance matrix of the data and provide the two key pieces of information needed to perform the dimensionality reduction: the **direction** and the **magnitude** of the new feature axes.

### Explanation

#### 1. Eigenvectors (The "Direction" Vectors)

- **Definition:** An eigenvector of a matrix (in PCA, the covariance matrix) is a non-zero vector that, when multiplied by the matrix, results in a scaled version of the original vector. Its direction does not change.
- **Role in PCA:** The eigenvectors of the covariance matrix define the **directions of the principal components**. They are orthogonal (perpendicular) to each other and represent the new axes of the transformed feature space. The first eigenvector points in the direction of the maximum variance in the data, the second points in the direction of the second-most variance (orthogonal to the first), and so on.

2.

#### 3. Eigenvalues (The "Magnitude" Scalars)

- **Definition:** An eigenvalue is the scalar factor by which the eigenvector is scaled after being multiplied by the matrix.
- **Role in PCA:** Each eigenvalue corresponds to a specific eigenvector. It represents the **amount of variance** in the data that is captured by that eigenvector (principal component). A large eigenvalue means that its corresponding principal component explains a large amount of the total variance in the original dataset.

4.

### Step-by-step Explanation of the Logic

1. After computing the covariance matrix of the data, we solve the equation  $\text{Cov}(X) * v = \lambda * v$ , where  $v$  is an eigenvector and  $\lambda$  (lambda) is its corresponding eigenvalue.
  2. The solution yields a set of eigenvector-eigenvalue pairs.
  3. We sort these pairs in descending order based on the eigenvalues.
  4. The eigenvector with the highest eigenvalue is **Principal Component 1 (PC1)**. The eigenvector with the second-highest eigenvalue is **Principal Component 2 (PC2)**, and so on.
  5. By selecting the top  $k$  eigenvectors, we choose the  $k$  most significant dimensions that capture the most information about the data.
- 

## Question 3

**Describe the role of the covariance matrix in PCA.**

### Theory

The covariance matrix is a fundamental building block in PCA. It is a square matrix that summarizes the **variance** of each feature and the **covariance** between each pair of features in the dataset. Its role is to provide the basis from which the principal components—the directions of maximum variance—are computed.

### Explanation

1. **What it Represents:**
  - **Diagonal Elements:** The elements on the main diagonal of the covariance matrix represent the **variance** of each individual feature. A high value indicates that the feature has a large spread.
  - **Off-Diagonal Elements:** The off-diagonal elements represent the **covariance** between pairs of features.
    - A positive covariance indicates that two features tend to increase or decrease together.
    - A negative covariance indicates that as one feature increases, the other tends to decrease.
    - A covariance near zero suggests that the two features are linearly independent.
  -
- 2.
3. **Why it is Central to PCA:**
  - **Capturing Relationships:** PCA's goal is to find new, uncorrelated axes that capture the maximum variance. The covariance matrix perfectly encapsulates the information needed for this task: it tells us how spread out the data is (variance) and in which directions the data varies together (covariance).

- **Source of Eigenvectors/Eigenvalues:** The principal components are found by performing an **eigen-decomposition** of the covariance matrix. The eigenvectors of the covariance matrix are the principal components, and the eigenvalues represent the variance explained by these components.
- In essence, the covariance matrix defines the shape and orientation of the data's multidimensional "cloud" of points. PCA then finds the principal axes of this cloud.

4.

## Best Practices

- **Data Centering:** Before computing the covariance matrix, the data must be **centered** by subtracting the mean from each feature. This ensures that the resulting principal components represent directions of variance around the data's center of mass.
  - **Data Scaling:** Because the covariance matrix is sensitive to the scale of the features, the data should also be **standardized** (to unit variance) before computing the covariance matrix. If this is done, the covariance matrix becomes equivalent to the correlation matrix.
- 

## Question 4

### What is the variance explained by a principal component?

#### Theory

The variance explained by a principal component is a measure of how much of the total variability in the original dataset is captured by that single component. It is directly represented by the **eigenvalue** corresponding to that principal component. The sum of the variances explained by all principal components is equal to the total variance of the original dataset.

#### Explanation

1. **Calculation:**
  - First, calculate the eigenvalues of the data's covariance matrix.
  - The total variance of the dataset is the sum of all these eigenvalues.
  - The variance explained by a single principal component  $i$  is simply its eigenvalue,  $\lambda_i$ .
- 2.
3. **Explained Variance Ratio:**
  - It is often more intuitive to express this as a percentage or ratio. The **explained variance ratio** for a principal component  $i$  is calculated as:  

$$\text{Explained Variance Ratio (PC}_i\text{)} = \lambda_i / (\sum \lambda_j \text{ for all } j)$$

- This ratio tells us the proportion of the total information (variance) that is contained within that specific component.
- 4.
5. **Cumulative Explained Variance:**
- To decide how many components to keep, we often look at the **cumulative explained variance**. This is the sum of the explained variance ratios of the top k principal components.
  - For example, if the first three components have a cumulative explained variance of 0.95, it means that by keeping only these three components, we have successfully retained 95% of the total variance (information) of the original high-dimensional dataset.
- 6.

## Use Cases

- **Dimensionality Reduction:** The explained variance ratio is the primary tool for deciding how many principal components to retain. A common practice is to keep enough components to explain a desired amount of variance (e.g., 90%, 95%, or 99%).
  - **Scree Plot:** This is a plot of the eigenvalues (explained variances) in descending order. It helps visualize the relative importance of each component and can be used to identify an "elbow" point, suggesting a natural cut-off for the number of components to keep.
- 

## Question 5

### How does scaling of features affect PCA?

#### Theory

Feature scaling is a **critical preprocessing step** for PCA. Because PCA is a variance-maximizing algorithm, it is highly sensitive to the scale of the input features. If features are on different scales, the one with the largest variance will dominate the first principal component, leading to biased and misleading results.

#### Explanation

1. **The Problem of Unequal Scales:**
- PCA identifies principal components based on the directions of maximum variance in the data.
  - Consider a dataset with two features: age (in years, e.g., 20-70) and income (in dollars, e.g., 30,000-150,000). The variance of income will be orders of magnitude larger than the variance of age.
  - When PCA is applied to this unscaled data, it will find that the direction of maximum variance is almost entirely aligned with the income axis. The first

principal component will be almost entirely determined by income, and the information from the age feature will be largely ignored.

2.

3. **The Solution: Standardization:**

- To prevent this, we must **standardize** the data before applying PCA. Standardization (or Z-score normalization) rescales each feature to have a mean of 0 and a standard deviation of 1.
- By doing this, all features are put on a comparable scale. This ensures that each feature has an equal opportunity to contribute to the principal components, and the resulting components will reflect the true underlying correlations and variance structure of the data, rather than being artifacts of arbitrary measurement units.

4.

### Best Practices

- **Always Standardize:** Unless you have a very specific reason not to (e.g., all features are already in the same units and you intentionally want to prioritize variance from certain features), you should always use a StandardScaler from a library like scikit-learn to preprocess your data before feeding it into a PCA model.
  - **Covariance vs. Correlation:** Applying PCA on standardized data is mathematically equivalent to performing the eigen-decomposition on the **correlation matrix** instead of the covariance matrix. The correlation matrix is essentially the covariance matrix of standardized variables.
- 

## Question 6

### What is the difference between PCA and Factor Analysis?

#### Theory

PCA and Factor Analysis (FA) are both linear dimensionality reduction techniques, but they have different underlying assumptions and goals. PCA is a purely mathematical transformation of the data, while FA is a statistical model that assumes an underlying latent variable structure.

#### Key Differences

Feature	Principal Component Analysis (PCA)	Factor Analysis (FA)
---------	---------------------------------------	----------------------

<b>Primary Goal</b>	To account for the <b>total variance</b> in the observed variables. It is a data simplification tool.	To explain the <b>common variance</b> (covariance) among observed variables via a smaller number of unobserved <b>latent factors</b> . It is an explanatory model.
<b>Underlying Model</b>	<b>No formal model.</b> It is an algorithmic transformation. The principal components are linear combinations of the observed variables.	<b>Assumes a latent variable model:</b> Observed Variables = $f(\text{Latent Factors}) + \text{Error}$ . It assumes the observed correlations are due to these hidden factors.
<b>Variance Decomp.</b>	Decomposes the total variance.	Decomposes the variance into <b>common variance</b> (shared by factors) and <b>unique variance</b> (error + variable-specific).
<b>Direction of Fit</b>	The principal components are defined by the observed variables. $\text{PC} = w_1X_1 + w_2X_2 + \dots$	The observed variables are defined by the latent factors. $X_1 = b_1F_1 + b_2F_2 + \dots + e$
<b>Output</b>	Principal components are orthogonal and represent directions of maximum variance.	The factors are estimated and are not necessarily orthogonal. The output also includes "factor loadings," which are correlations between variables and factors.

## Use Cases

- **Use PCA when:**
  - The primary goal is **dimensionality reduction**.
  - You want to create a set of uncorrelated features for a subsequent machine learning model.
  - You want a simple, deterministic way to summarize the data.
- 
- **Use Factor Analysis when:**
  - You want to understand the **underlying latent structure** of your data.
  - You are developing a psychological scale or survey and want to identify underlying constructs (e.g., identifying "intelligence" and "creativity" as latent factors from test scores).
- 

**Analogy:** If your variables are `test_score_math`, `test_score_logic`, `test_score_reading`, and `test_score_vocabulary`:

- **PCA** would just find the best way to combine these into fewer components to summarize the data. PC1 might be a weighted average of all four.
  - **FA** would try to find if there are underlying latent factors like "Quantitative Ability" (explaining math/logic scores) and "Verbal Ability" (explaining reading/vocab scores).
- 

## Question 7

**Can you explain the Singular Value Decomposition (SVD) and its relationship with PCA?**

### Theory

**Singular Value Decomposition (SVD)** is a fundamental matrix factorization technique in linear algebra. It asserts that any  $m \times n$  matrix  $X$  can be decomposed into the product of three other matrices:  $X = U\Sigma V^T$ .

- $U$ : An  $m \times m$  orthogonal matrix whose columns are the left-singular vectors.
- $\Sigma$  (Sigma): An  $m \times n$  diagonal matrix containing the singular values of  $X$ .
- $V^T$ : The transpose of an  $n \times n$  orthogonal matrix  $V$  whose columns are the right-singular vectors.

The relationship between SVD and PCA is very direct: **SVD is a numerically stable and efficient method for calculating the principal components of a dataset.**

### The Relationship Explained

While PCA can be defined conceptually through the eigen-decomposition of the covariance matrix, in practice, most software implementations (including scikit-learn) use SVD because it is more robust.

Here's the connection for a centered data matrix  $X$ :

1. **Covariance Matrix:** The covariance matrix of  $X$  is  $(1/(n-1)) * X^T X$ .
2. **PCA's Goal:** PCA finds the eigenvectors of this covariance matrix.
3. **SVD's Role:** Let the SVD of the centered data matrix  $X$  be  $X = U\Sigma V^T$ .
  - The columns of the matrix  $V$  are the **eigenvectors** of  $X^T X$ . These are precisely the **principal components** (the directions).
  - The singular values in  $\Sigma$  are related to the eigenvalues of the covariance matrix by  $\lambda = \sigma^2 / (n-1)$ . The squared singular values are proportional to the eigenvalues and thus represent the **variance explained** by each principal component.
- 4.

### Why Use SVD for PCA?

- **Numerical Stability:** Computing the covariance matrix  $X^T X$  can lead to a loss of numerical precision, especially if the data has a high condition number. SVD works directly on the data matrix  $X$ , avoiding this potentially problematic step and providing more accurate results.
- **Efficiency:** For datasets where the number of features is much larger than the number of samples, SVD can be computationally faster than forming the full covariance matrix.
- **Unified Approach:** SVD provides a direct and elegant way to obtain the principal components, the projected data, and the explained variance all from a single decomposition.

In summary, SVD is the engine that powers modern implementations of PCA.

---

## Question 8

**What is meant by 'loading' in the context of PCA?**

### Theory

In PCA, **loadings** (or component loadings) are the **correlations between the original variables and the principal components**. They represent how much each original feature contributes to the creation of a particular principal component. A high loading value (close to +1 or -1) indicates that the feature has a strong relationship with that component.

### Explanation

1. **Interpretation:** Loadings help us understand the meaning of the principal components. After performing PCA, we get abstract components (PC1, PC2, etc.). By examining the loadings, we can interpret what these components represent in terms of our original features.
  - **Example:** Suppose we perform PCA on a dataset of car features. We might find that PC1 has high positive loadings for horsepower, engine\_size, and price, and high negative loadings for miles\_per\_gallon. We could then interpret PC1 as a "Performance/Power" component.
- 2.
3. **Calculation:** The loading of a variable  $j$  on a principal component  $i$  is the correlation between that variable and the component scores. It can also be calculated from the eigenvector  $v_i$  and the standard deviation  $s_j$  of the original variable:  

$$\text{Loading}(j, i) = v_i(j) * \sqrt{\lambda_i} / s_j$$
 For standardized data (which is the best practice), the loadings are simply the eigenvectors themselves.
4. **Loadings vs. Eigenvectors:**

- **Eigenvectors** define the mathematical direction of the principal components. Their values are the weights used to linearly combine the original variables to get the component scores.
  - **Loadings** provide a more interpretable, correlation-based view of the relationship between the original variables and the components.
- 5.

## Use Cases

- **Component Interpretation:** This is the primary use. Loadings are essential for assigning business or scientific meaning to the otherwise abstract principal components.
  - **Variable Importance:** Variables with high loadings on the first few important components are generally considered more significant in defining the data's variance structure.
- 

## Question 9

**Explain the process of eigenvalue decomposition in PCA.**

### Theory

Eigenvalue decomposition, or eigen-decomposition, is the mathematical process at the heart of PCA. It involves breaking down a square matrix (specifically, the covariance matrix of the data) into its constituent eigenvectors and eigenvalues. This decomposition reveals the fundamental structure of the linear transformation represented by the matrix, which in PCA corresponds to the directions and magnitudes of maximum variance.

### Step-by-Step Process

1. **Start with the Covariance Matrix (C):** After standardizing the data, the first step is to compute the  $p \times p$  covariance matrix  $C$  for the  $p$  features. This matrix is symmetric.
2. **The Eigenvalue Equation:** The goal is to find the eigenvectors  $v$  and eigenvalues  $\lambda$  that satisfy the characteristic equation:

$$C * v = \lambda * v$$

- $C$ : The covariance matrix.
- $v$ : An eigenvector (a  $p \times 1$  column vector).
- $\lambda$ : A corresponding eigenvalue (a scalar).

This equation states that when the covariance matrix  $C$  acts on an eigenvector  $v$ , it only scales  $v$  by its eigenvalue  $\lambda$ , without changing its direction.

3.

4. **Finding the Eigenvalues:** To solve this, we rearrange the equation:

$$(C - \lambda I) * v = 0$$

where  $I$  is the identity matrix. For this equation to have a non-zero solution for  $v$ , the

matrix  $(C - \lambda I)$  must be singular, which means its determinant must be zero:

$$\det(C - \lambda I) = 0$$

Solving this determinant equation (the "characteristic polynomial") gives us the p eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_p$ .

5. **Finding the Eigenvectors:** For each calculated eigenvalue  $\lambda_i$ , we plug it back into the equation  $(C - \lambda_i I) * v = 0$  and solve for the corresponding eigenvector  $v_i$ .
6. **Result of Decomposition:** The final result is a set of p eigenvector-eigenvalue pairs. The eigenvectors  $(v_1, v_2, \dots, v_p)$  are mutually orthogonal and represent the principal components. The eigenvalues  $(\lambda_1, \lambda_2, \dots, \lambda_p)$  represent the variance captured by each of these components.

This decomposition provides all the necessary information to transform the original data into the new principal component space.

---

## Question 10

**What are the limitations of PCA when it comes to handling non-linear relationships?**

### Theory

The primary limitation of standard PCA is that it is a **linear** algorithm. It assumes that the underlying structure of the data can be captured by a linear subspace. Consequently, PCA is incapable of discovering or effectively representing **non-linear relationships** among variables.

### Explanation of the Limitation

1. **Linear Projections:** PCA works by projecting the data onto a set of orthogonal linear axes (the principal components). If the data lies on or near a non-linear manifold (e.g., a curve like a "Swiss roll" or an "S" shape), any linear projection will fail to capture this structure. Projecting a Swiss roll onto a 2D plane using PCA would just result in a flattened, overlapping blob of points, completely losing the original coiled structure.
2. **Maximizing Variance Linearly:** The objective of PCA is to find linear combinations of features that maximize variance. This works well if the primary sources of variation in the data are linear. However, if the interesting patterns are non-linear, PCA will miss them.
3. **Example Scenario:** Consider a dataset where points are arranged in a circle. There are no linear directions of high variance. PCA would likely produce two components with roughly equal variance, and the projection would not reveal the circular pattern.

### Solutions and Alternatives

When non-linear relationships are suspected, standard PCA is not the right tool. The following alternatives are more appropriate:

- **Kernel PCA (kPCA):** This is the most direct extension of PCA for non-linear data. It uses the "kernel trick" to implicitly map the data into a much higher-dimensional feature space where the relationships might become linear. PCA is then performed in this higher-dimensional space, allowing it to capture non-linear structures in the original data.
  - **Manifold Learning Algorithms:** These techniques are specifically designed to find the underlying low-dimensional non-linear manifold on which the data lies.
    - **t-SNE (t-Distributed Stochastic Neighbor Embedding):** Excellent for visualization, as it focuses on preserving local neighborhood structures.
    - **Isomap:** Preserves the geodesic distances between points on the manifold.
  - 
  - **Autoencoders:** A type of neural network that can learn powerful non-linear embeddings of data. The encoder part of a trained autoencoder can be used as a non-linear dimensionality reduction tool.
- 

## Question 11

**Explain the curse of dimensionality and how PCA can help to mitigate it.**

### Theory

The **curse of dimensionality** refers to a collection of problems that arise when analyzing and modeling data in high-dimensional spaces. As the number of features (dimensions) increases, the data becomes increasingly sparse, distance metrics become less meaningful, and the performance of many machine learning algorithms degrades.

### Problems Caused by the Curse of Dimensionality

1. **Data Sparsity:** The volume of the feature space grows exponentially with the number of dimensions. A fixed number of data points becomes increasingly sparse, making it difficult to find meaningful patterns or clusters.
2. **Distance Concentration:** In high dimensions, the distances between most pairs of points tend to become almost equal. This makes distance-based algorithms like K-means and k-NN less effective, as the concept of a "nearest" neighbor becomes ill-defined.
3. **Overfitting:** With more features than samples, many supervised learning models are prone to overfitting. They learn the noise in the data rather than the true underlying signal.
4. **Computational Cost:** The computational resources required to process and model the data increase significantly with the number of dimensions.

### How PCA Mitigates the Curse

PCA helps to mitigate these issues by reducing the number of dimensions while retaining the most important information.

1. **Dimensionality Reduction:** By projecting the data onto a lower-dimensional subspace defined by the top k principal components, PCA directly reduces the number of features. This addresses the computational cost and sparsity issues.
2. **Noise Reduction:** Principal components are ordered by the amount of variance they explain. The components with low variance often correspond to noise in the data. By discarding these later components, PCA acts as a denoising filter, which can improve the performance of subsequent models.
3. **Combating Overfitting:** By reducing the number of features, PCA simplifies the problem for a subsequent supervised learning model. With fewer features, the model is less likely to overfit the training data.
4. **Improving Model Performance:** By transforming correlated features into a smaller set of uncorrelated components, PCA can improve the stability and performance of models that are sensitive to multicollinearity (like linear regression).

In essence, PCA provides a new, lower-dimensional representation of the data that is more compact, less noisy, and more suitable for analysis and modeling, thereby breaking the curse.

---

## Question 12

**How does PCA handle missing values in the data?**

### Theory

Standard implementations of PCA **cannot handle missing values**. The mathematical operations at its core—calculating the covariance matrix and performing SVD or eigen-decomposition—require a complete numerical matrix. If missing values (e.g., NaNs) are present in the dataset, the algorithm will fail.

### Strategies to Address Missing Values Before PCA

The responsibility for handling missing values falls on the user during the preprocessing stage.

1. **Deletion:**
  - **Method:** Remove rows or columns containing missing values.
  - **Pitfall:** This is a simple but often poor solution, as it can lead to significant data loss and introduce bias if the data is not missing completely at random.
- 2.
3. **Imputation (The Standard Approach):**

- **Method:** Replace missing values with an estimated value. This must be done **before** applying PCA.
  - **Simple Imputation:** Use the mean, median, or mode of the respective feature. Median is often preferred for numerical data as it is robust to outliers.
  - **Advanced Imputation:** Use more sophisticated methods like k-NN imputation or MICE for more accurate estimates.
  - **Caution:** Imputation can artificially reduce the natural variance of the data and weaken correlations, which can affect the results of PCA.
- 4.
5. **Specialized PCA Variants (Probabilistic PCA):**
- **Concept:** More advanced versions of PCA exist that can handle missing values directly as part of their model.
  - **Probabilistic PCA (PPCA):** This is a probabilistic formulation of PCA that uses an Expectation-Maximization (EM) algorithm. It can gracefully handle missing data by marginalizing over the missing values during the E-step and M-step, providing a more principled approach than simple imputation.
  - **Implementation:** While not the default, some libraries offer implementations of these advanced variants.
- 6.

### **Best Practice**

For most applications, the recommended workflow is:

1. Split the data into training and test sets.
  2. Learn an imputer (e.g., SimpleImputer) on the **training data only**.
  3. Transform both the training and test data using the fitted imputer.
  4. Apply PCA on the imputed training data.
  5. Transform the imputed test data using the fitted PCA model.
- 

## **Question 13**

**What is the difference between PCA and t-SNE for dimensionality reduction?**

### **Theory**

PCA and t-SNE (t-Distributed Stochastic Neighbor Embedding) are both dimensionality reduction techniques, but they are fundamentally different in their goals, methods, and appropriate use cases. PCA is a linear technique for data projection, while t-SNE is a non-linear technique primarily for data visualization.

### **Key Differences**

Feature	Principal Component Analysis (PCA)	t-SNE (t-Distributed Stochastic Neighbor Embedding)
<b>Primary Goal</b>	<b>Preserve global variance.</b> It finds directions of maximum variance to create a lower-dimensional representation of the data.	<b>Preserve local structure.</b> It models the similarity between neighboring points to create a low-dimensional embedding for visualization.
<b>Linearity</b>	<b>Linear.</b> It uses a linear transformation (projection) of the data.	<b>Non-linear.</b> It can capture complex, non-linear manifold structures.
<b>Output Interpretation</b>	The distances and angles between points in the PCA plot are meaningful and relate to the overall data structure. The axes (principal components) are also interpretable.	The output is purely for <b>visualization</b> . Distances between clusters might be meaningful, but the size of clusters and the distances between points within a cluster are not. The axes are not interpretable.
<b>Determinism</b>	<b>Deterministic.</b> Always produces the same output for the same data.	<b>Stochastic.</b> Results can vary slightly between runs due to random initialization.
<b>Computational Cost</b>	Relatively fast and scalable.	Computationally intensive ( $O(n \log n)$ or $O(n^2)$ ), not suitable for very large datasets.
<b>Primary Use Case</b>	<b>Dimensionality reduction</b> for machine learning models (feature engineering, noise reduction).	<b>High-dimensional data visualization</b> (e.g., visualizing the output of a neural network layer).

## Best Practices

- **Do not use t-SNE for clustering:** Because t-SNE can create the visual appearance of clusters where none exist, it is not a reliable tool for clustering itself.
- **Use PCA before t-SNE:** For very high-dimensional data, it is a common and effective practice to first use PCA to reduce the dimensionality to a moderate level (e.g., 50 dimensions) and then use t-SNE to reduce it further to 2 or 3 dimensions for visualization. This speeds up t-SNE and can help it find a better global structure.

## Question 14

**Explain how PCA can be used as a noise reduction technique.**

## Theory

PCA can be an effective technique for noise reduction based on the assumption that the signal (the important information) is concentrated in the directions of high variance, while the noise is distributed across all dimensions and contributes more to the directions of low variance. By reconstructing the data using only the principal components with the highest variance, PCA can filter out this noise.

## Explanation

1. **The Underlying Assumption:** The core idea is the **signal-to-noise ratio**. In many real-world datasets, the true underlying signal is responsible for the largest variations in the data. Random noise, on the other hand, adds small, random perturbations to all features, contributing primarily to the components with low variance.
2. **The Filtering Process:**
  - o **Decomposition:** PCA is first performed on the noisy data. It decomposes the data into a set of principal components, ordered by their explained variance (eigenvalues).
  - o **Component Selection:** The components with high eigenvalues (high variance) are assumed to represent the signal. The components with low eigenvalues (low variance) are assumed to represent the noise.
  - o **Reconstruction:** The data is then reconstructed using **only the top k principal components** that correspond to the signal. The information from the later, noisy components is discarded. The reconstructed data will be a "denoised" version of the original data, lying in the lower-dimensional subspace defined by the signal components.
- 3.

## Code Example (Conceptual)

```
code Python
downloadcontent_copyexpand_less
from sklearn.decomposition import PCA
import numpy as np

# Assume X_noisy is your noisy dataset
# 1. Fit PCA and decide how many components represent the signal
# For example, keep components that explain 95% of the variance
pca = PCA(n_components=0.95)
X_transformed = pca.fit_transform(X_noisy)

# 2. Reconstruct the data from the selected components
# This projects the data back into the original feature space, but without the noise.
X_reconstructed_denoised = pca.inverse_transform(X_transformed)
```

## Pitfalls and Limitations

- **The Assumption May Be Wrong:** This technique relies heavily on the assumption that low variance equals noise. In some cases, important but subtle information may be contained in the lower-variance components. Discarding them could be detrimental.
  - **Information Loss:** Denoising via PCA is inherently a lossy process. The reconstructed data is not identical to the original "clean" data, but an approximation of it.
- 

## Question 15

Describe how you would apply PCA for visualization purposes.

### Theory

PCA is a powerful and widely used technique for visualizing high-dimensional data. The human eye can only perceive 2D or 3D space. PCA allows us to create a lower-dimensional "shadow" or projection of a high-dimensional dataset that preserves the maximum possible amount of its variance, making it suitable for plotting and visual inspection.

### Step-by-Step Procedure

1. **Select the Data:** Identify the high-dimensional dataset you want to visualize. For example, the Iris dataset has 4 features, or a customer dataset might have dozens of features.
2. **Standardize the Features:** This is a crucial first step. Use a StandardScaler to ensure all features are on the same scale, so that none dominate the PCA calculation.
3. **Apply PCA:**
  - Instantiate a PCA model from a library like scikit-learn.
  - Set the number of components you want to reduce to. For visualization, this will be `n_components=2` for a 2D plot or `n_components=3` for a 3D plot.
  - Fit the PCA model to the standardized data and then transform the data to get the principal component scores.
- 4.
5. **Create the Visualization:**
  - The result of the transformation will be a new matrix where each data point is now represented by 2 (or 3) coordinates instead of its original high number of features.
  - Create a **scatter plot** using these new coordinates. The x-axis will be the first principal component (PC1), and the y-axis will be the second principal component (PC2).

- **Enhance with Labels:** To make the plot informative, color the points according to a known label, such as the species in the Iris dataset or a customer segment. This helps to visually assess whether the original classes are well-separated in the principal component space.

6.

### Code Example (Conceptual)

```
code Python
downloadcontent_copyexpand_less
IGNORE_WHEN COPYING_START
IGNORE_WHEN COPYING_END
    from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler
from sklearn.datasets import load_iris
import matplotlib.pyplot as plt

# 1. Load data
iris = load_iris()
X = iris.data
y = iris.target

# 2. Standardize
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# 3. Apply PCA for 2D visualization
pca = PCA(n_components=2)
X_pca = pca.fit_transform(X_scaled)

# 4. Create scatter plot
plt.figure(figsize=(8, 6))
scatter = plt.scatter(X_pca[:, 0], X_pca[:, 1], c=y, cmap='viridis')
plt.xlabel('First Principal Component (PC1)')
plt.ylabel('Second Principal Component (PC2)')
plt.title('PCA of Iris Dataset')
plt.legend(handles=scatter.legend_elements()[0], labels=iris.target_names)
plt.show()
```

---

This plot would show how the different Iris species form distinct clusters in the 2D space defined by the first two principal components.

# Question 16

What are the advantages and drawbacks of kernel PCA compared to linear PCA?

## Theory

**Kernel PCA (kPCA)** is an extension of standard PCA that uses the "kernel trick" to perform dimensionality reduction on non-linear data. It offers the significant advantage of capturing complex, non-linear structures but comes with its own set of trade-offs.

### Advantages of Kernel PCA

1. **Handles Non-Linearity:** This is its primary advantage. By using kernels (like the Radial Basis Function (RBF) or polynomial kernels), kPCA implicitly maps the data into a very high-dimensional feature space where non-linear relationships in the original space may become linear. It can then separate these non-linear structures, which linear PCA would fail to do.
  - o **Use Case:** Ideal for datasets with complex geometries like the "Swiss roll," concentric circles, or intertwined spirals.
- 2.
3. **More Powerful Representation:** kPCA can learn a more powerful and descriptive low-dimensional representation for data that does not conform to the linear assumptions of standard PCA.

### Drawbacks of Kernel PCA

1. **Higher Computational Cost:** kPCA is computationally more expensive than linear PCA. It requires the computation of an  $n \times n$  kernel matrix, which can be prohibitively large for datasets with a large number of samples ( $n$ ). Its complexity is typically  $O(n^2)$  or  $O(n^3)$ , compared to the more scalable standard PCA.
2. **Parameter Tuning:** kPCA introduces new hyperparameters that need to be tuned, such as the choice of the kernel (e.g., RBF, polynomial, sigmoid) and the kernel's specific parameters (e.g., the gamma parameter for the RBF kernel). The performance of kPCA is highly sensitive to these choices.
3. **Loss of Interpretability:** The components derived from kPCA do not have a clear interpretation in terms of the original features. Unlike linear PCA where we can use loadings to understand the components, the components in kPCA exist in an implicit high-dimensional space and are harder to relate back to the original variables.
4. **No Direct Reconstruction:** Standard kPCA does not provide a direct way to reconstruct the original data from the projected components (the "pre-image problem"), although approximate solutions exist.

**Conclusion:** Use **linear PCA** as a fast, scalable, and interpretable baseline. Use **Kernel PCA** when you have strong reason to believe the data contains important non-linear structures and you can afford the higher computational cost and parameter tuning effort.

---

## Question 17

**Explain how you would apply PCA in a stock market data analysis situation.**

### Theory

In stock market analysis, PCA can be a powerful tool to reduce the dimensionality of data containing the returns of many different stocks. Stock prices are often highly correlated (e.g., stocks in the same sector tend to move together). PCA can distill this complex, correlated data into a few key "factors" or components that drive the overall market movement.

### Step-by-Step Application

1. **Data Preparation:**

- **Select Stocks:** Choose a portfolio of stocks to analyze (e.g., all stocks in the S&P 500).
- **Calculate Returns:** Convert the raw stock prices into daily or weekly returns. Returns are more stationary and suitable for analysis than prices. The data matrix  $X$  would have  $n$  days as rows and  $p$  stocks as columns.
- **Standardize:** Standardize the returns for each stock to have a mean of 0 and a standard deviation of 1. This ensures that highly volatile stocks do not dominate the analysis.

2.

3. **Apply PCA:**

- Perform PCA on the standardized returns matrix.
- The result will be a set of principal components, each representing an orthogonal source of variation in the stock market returns.

4.

5. **Interpretation of Principal Components:**

- **PC1 (The Market Component):** The first principal component will almost always represent the **overall market movement**. It will have positive loadings on nearly all stocks and will capture the largest portion of the total variance. A high score on PC1 on a given day means the market as a whole went up.
- **PC2, PC3, etc. (Sector/Style Components):** The subsequent components represent more nuanced factors that are independent of the overall market. By examining the **loadings**, we can interpret their meaning.
  - For example, PC2 might have high positive loadings on technology stocks and high negative loadings on utility stocks. This could be interpreted as a "Growth vs. Value" or "Tech vs. Defensive" factor. A high score on PC2 would mean tech stocks outperformed utilities on that day.
  - Other components might represent factors like company size, interest rate sensitivity, or commodity prices.

- - 6.

## Use Cases in Finance

- **Factor-Based Investing:** The principal components can be used as the basis for creating factor-based investment strategies.
  - **Risk Management:** By understanding the key factors driving portfolio returns, risk managers can better hedge against specific risks (e.g., exposure to the "market factor" or a "sector factor").
  - **Portfolio Diversification:** PCA can reveal hidden correlations and help construct a portfolio that is truly diversified across different sources of risk.
- 

## Question 18

**Describe a scenario where using PCA might be detrimental to the performance of a machine learning model.**

### Theory

While PCA is often beneficial, its application can be detrimental if its core assumption—that high variance corresponds to high signal—is violated. If the features that are most important for predicting the target variable have low variance, PCA might discard them, leading to a significant drop in model performance.

### Scenario: Low-Variance, High-Impact Features

Imagine a classification problem where you are trying to predict whether a customer will **churn** (cancel their subscription). The dataset has many features, including:

- total\_spend: Total amount spent by the customer (high variance, range: \$10 - \$10,000).
- monthly\_data\_usage: Data used per month (high variance, range: 1GB - 500GB).
- complaints\_filed: A binary feature (0 or 1) indicating if the customer has ever filed a formal complaint (very low variance).

### The Problem:

- The vast majority of customers have never filed a complaint, so the complaints\_filed feature has very low variance.
- However, having filed a complaint is a **very strong predictor of churn**. It is a low-variance feature with high predictive power.

### Applying PCA:

1. When PCA is applied to this dataset, it will focus on capturing the directions of high variance, which will be dominated by total\_spend and monthly\_data\_usage.
2. The complaints\_filed feature, due to its low variance, will contribute very little to the first few principal components.
3. If we then reduce the dimensionality by keeping only the top components (which is the whole point of using PCA), we will effectively **discard the information** from the complaints\_filed feature.

#### **The Detrimental Outcome:**

- A machine learning model (e.g., Logistic Regression or a Random Forest) trained on the PCA-transformed data will likely have **worse performance** than a model trained on the original data.
- The model will be deprived of one of its most powerful predictors because PCA, being an unsupervised technique, has no knowledge of the target variable (churn) and incorrectly equated low variance with low importance.

**Conclusion:** PCA should be used with caution when the link between variance and predictive importance is not guaranteed. In supervised learning contexts, feature selection methods that explicitly consider the target variable (like feature importance from a Random Forest or recursive feature elimination) are often safer and more effective.

---

## **Question 19**

### **What are the best practices in visualizing the components obtained from PCA?**

#### **Theory**

Visualizing the results of PCA is crucial for understanding the structure of the data and the meaning of the components. Best practices focus on creating plots that are not only accurate but also highly interpretable.

#### **Best Practices for Visualization**

##### **1. 2D Scatter Plot of Components:**

- **Method:** The most common visualization is a scatter plot of the data points projected onto the first two principal components (PC1 vs. PC2).
- **Best Practice:**
  - **Label Axes Clearly:** Label the x-axis as "First Principal Component (PC1)" and the y-axis as "Second Principal Component (PC2)".
  - **Include Explained Variance:** Add the percentage of variance explained by each component to the axis labels (e.g., "PC1 (Explains 45.2% of

variance)"). This gives the viewer context about how much information the plot represents.

- **Color-Code by a Target Variable:** If you have external labels (e.g., class labels, cluster assignments), color the points accordingly. This is the most powerful way to see if the components effectively separate the groups.

○

2.

3. **Scree Plot:**

- **Method:** A bar or line plot showing the explained variance (or explained variance ratio) for each principal component, in descending order.
- **Best Practice:** This plot is essential for deciding how many components to keep. Look for the "elbow" point where the explained variance drops off sharply. It's also common to overlay a line showing the cumulative explained variance.

4.

5. **Biplot:**

- **Method:** A biplot is a more advanced visualization that **overlays** the component scores (the scatter plot of the data points) with the component loadings (vectors representing the original features).
- **Interpretation:**
  - The data points show the structure of the samples.
  - The vectors (arrows) show how the original features contribute to the components. A long vector indicates a feature has a strong influence. The direction of the vector shows its correlation with the components.
- 
- **Best Practice:** A biplot provides a rich, single view of both the samples and the features, but it can become cluttered with many features.

6.

7. **Loading Plot:**

- **Method:** A bar chart or scatter plot showing the loading values of each original feature for a given principal component.
- **Best Practice:** This is used specifically for interpreting the meaning of the components. For example, a bar chart of the loadings for PC1 would clearly show which features are most influential in its definition.

8.

By combining these different visualizations, a data scientist can gain a comprehensive understanding of what PCA has revealed about the dataset.

---

## Question 20

**Explain how Incremental PCA differs from standard PCA and when you would use it.**

## Theory

**Incremental PCA (IPCA)** is a variant of PCA designed to handle datasets that are too large to fit into a single machine's memory. Unlike standard PCA, which requires the entire dataset to be loaded at once to compute the covariance matrix or perform SVD, IPCA processes the data in mini-batches.

## How Incremental PCA Works

1. **Mini-Batch Processing:** IPCA breaks the large dataset into smaller chunks or mini-batches that can fit in memory.
2. **Iterative Updates:** It processes one mini-batch at a time. For each batch, it updates its estimate of the principal components and the explained variance. It does this without looking at the full data, only at the current mini-batch and the summary statistics from the previous batches.
3. **Approximation:** The final result is an approximation of what standard PCA would have produced. By using a reasonable batch size, the approximation is typically very close to the exact solution.

## Comparison with Standard PCA

Feature	Standard PCA	Incremental PCA (IPCA)
<b>Data Handling</b>	Requires the entire dataset to be in memory.	Processes data in mini-batches; does not require the full dataset in memory.
<b>Memory Usage</b>	High. Proportional to $n_{samples} * n_{features}$ .	Low. Proportional to $batch\_size * n_{features}$ .
<b>Computation</b>	Performs a single, exact SVD on the full dataset.	Performs partial SVD updates for each mini-batch.
<b>Result</b>	Exact.	An approximation of the exact solution. The quality depends on the batch size.
<b>Use Case</b>	For datasets that comfortably fit in RAM.	For <b>very large datasets</b> that cannot fit in memory or for <b>streaming data</b> applications where data arrives continuously.

## When to Use Incremental PCA

- **Big Data:** This is the primary use case. When your dataset is larger than your available RAM (e.g., tens or hundreds of gigabytes), IPCA is the only feasible option for performing PCA on a single machine.
- **Online Learning / Streaming Data:** IPCA can be used in an online learning setting. As new data arrives in a stream, it can be treated as a new mini-batch to continuously update the PCA model without retraining from scratch.

In scikit-learn, IncrementalPCA has an API very similar to PCA, using `partial_fit()` to process each mini-batch.

---

## Question 21

**How does Generalized PCA differ from standard PCA and what are its applications?**

### Theory

**Generalized PCA (GPCA)**, often associated with the Generalized Linear Model (GLM) framework, is an extension of standard PCA that can handle data that comes from distributions other than the Gaussian distribution. Standard PCA is implicitly designed for data that is normally distributed, as its objective of minimizing squared reconstruction error is equivalent to maximizing the likelihood under a Gaussian assumption.

GPCA extends this idea to other distributions in the **exponential family**, such as the Bernoulli, Poisson, and Gamma distributions.

### Key Differences from Standard PCA

1. **Data Distribution Assumption:**
  - **Standard PCA:** Assumes the data follows a Gaussian distribution. The loss function is the squared error.
  - **GPCA:** Can handle data from various distributions in the exponential family (e.g., binary data, count data). The loss function is derived from the log-likelihood of the assumed distribution (e.g., logistic loss for Bernoulli data, Poisson loss for count data).
- 2.
3. **Linearity:**
  - **Standard PCA:** Finds a linear subspace. The relationship between the observed data and the low-dimensional representation is linear.
  - **GPCA:** The relationship is defined by a **link function** from the GLM framework. For example, for binary data (Bernoulli distribution), a logit link function would be used. This allows GPCA to model non-linear relationships between the latent space and the observed data.
- 4.

### Applications

GPCA is particularly useful when the data is clearly non-Gaussian.

- **Binary Data:** For datasets where all features are binary (0/1), such as survey responses (yes/no) or user interaction data (clicked/not clicked). GPCA with a Bernoulli assumption (sometimes called Logistic PCA) would be more appropriate than standard PCA.
- **Count Data:** For datasets representing counts, such as word counts in documents (text analysis) or the number of events in a time interval. GPCA with a Poisson assumption can effectively reduce the dimensionality of this type of data.
- **Positive Continuous Data:** For data that is strictly positive and skewed, like financial amounts or reaction times, GPCA with a Gamma distribution assumption might be suitable.

In essence, GPCA provides a more statistically sound and flexible framework for dimensionality reduction when the assumptions of standard PCA are violated by the nature of the data.

---

## Question 22

**Define PCA and its optimization objective.**

### Theory

Principal Component Analysis (PCA) is an unsupervised linear transformation technique used for dimensionality reduction. It projects a dataset onto a lower-dimensional subspace while preserving as much of the original data's variability as possible.

This goal can be framed through two equivalent optimization objectives:

1. **Maximizing Variance:** Find the projection that maximizes the variance of the projected data.
2. **Minimizing Reconstruction Error:** Find the projection that minimizes the average squared distance between the original data points and their reconstructed approximations.

### Explanation of Optimization Objectives

#### 1. Variance Maximization Objective

- **Goal:** Find a set of orthogonal unit vectors (the principal components)  $v_1, v_2, \dots, v_d$  such that the variance of the data projected onto these vectors is maximized.
- **Procedure:**
  - The first principal component  $v_1$  is the unit vector that maximizes the variance of the projected data  $\text{Var}(Xv_1)$ .
  - The second principal component  $v_2$  is the unit vector, orthogonal to  $v_1$ , that maximizes the variance of the projected data  $\text{Var}(Xv_2)$ .
  - This continues for all subsequent components.
-

- **Result:** This objective leads directly to finding the eigenvectors of the data's covariance matrix, as these are the directions of maximum variance.
- 2.
3. **Reconstruction Error Minimization Objective**
- **Goal:** Find a k-dimensional subspace such that when the data is projected onto it and then projected back into the original space, the error (distance) between the original points and their reconstructed versions is minimized.
  - **Procedure:** Find a projection that minimizes the average squared Euclidean distance:  $\sum ||x_i - \hat{x}_i||^2$ , where  $x_i$  is an original data point and  $\hat{x}_i$  is its reconstruction from the low-dimensional subspace.
  - **Result:** It can be mathematically proven that the subspace that minimizes this reconstruction error is the very same one spanned by the eigenvectors corresponding to the largest eigenvalues of the covariance matrix.
- 4.

**Conclusion:** Both objectives, although conceptually different, lead to the exact same mathematical solution: the subspace defined by the top k eigenvectors of the data's covariance matrix. Maximizing variance is the more common and intuitive way to explain PCA.

---

## Question 23

Show link between PCA and eigen-decomposition of covariance.

### Theory

The fundamental link between PCA and the eigen-decomposition of the covariance matrix is that the **principal components are precisely the eigenvectors of the covariance matrix**. This connection arises directly from PCA's objective of finding the directions that maximize the variance of the projected data.

### Mathematical Derivation

Let  $X$  be the centered data matrix of size  $n \times p$  ( $n$  samples,  $p$  features). The covariance matrix is  $C = (1/(n-1)) * X^T X$ .

Let  $v$  be a unit vector representing a direction onto which we project the data. The projected data is  $Xv$ .

1. **Objective:** We want to find the direction  $v$  that maximizes the variance of the projected data  $\text{Var}(Xv)$ .

$$\begin{aligned} \text{Var}(Xv) &= (1/(n-1)) * (Xv)^T (Xv) \\ &= (1/(n-1)) * v^T X^T X v \\ &= v^T * [(1/(n-1)) * X^T X] * v \\ &= v^T C v \end{aligned}$$

2. **Optimization Problem:** Our goal is to maximize  $v^T C v$  subject to the constraint that  $v$  is a unit vector, i.e.,  $v^T v = 1$ .
  3. **Lagrange Multiplier:** We can solve this constrained optimization problem using a Lagrange multiplier  $\lambda$ . The Lagrangian function is:  

$$L(v, \lambda) = v^T C v - \lambda(v^T v - 1)$$
  4. **Find the Gradient:** To find the maximum, we take the derivative of  $L$  with respect to  $v$  and set it to zero:  

$$\partial L / \partial v = 2Cv - 2\lambda v = 0$$
  

$$Cv - \lambda v = 0$$
  

$$Cv = \lambda v$$
  5. **Conclusion:** This is the **eigenvalue equation**. It shows that the vector  $v$  that maximizes the variance must be an **eigenvector** of the covariance matrix  $C$ .
    - To find which eigenvector gives the *maximum* variance, we can substitute  $Cv = \lambda v$  back into our variance formula:  

$$\text{Variance} = v^T Cv = v^T(\lambda v) = \lambda(v^T v) = \lambda * 1 = \lambda$$
    - This shows that the variance along the direction of an eigenvector  $v$  is simply its corresponding **eigenvalue  $\lambda$** .
    - Therefore, to maximize the variance, we must choose the eigenvector with the **largest eigenvalue**. This becomes our first principal component. The second principal component is the eigenvector with the second-largest eigenvalue, and so on.
- 6.
- 

## Question 24

**Explain variance maximization vs. reconstruction error minimization.**

### Theory

Variance maximization and reconstruction error minimization are two different but mathematically equivalent perspectives for understanding the optimization objective of PCA. Both frameworks lead to the same solution: the principal components are the top  $k$  eigenvectors of the data's covariance matrix.

### Variance Maximization View

- **Concept:** This view focuses on finding a low-dimensional representation that **preserves the most information** from the original data. In PCA, "information" is quantified as variance.
- **Objective:** Find a set of orthogonal axes (principal components) such that the variance of the data projected onto these axes is maximized.

- **Intuition:** We are looking for the directions in which the data is most "spread out." These directions are assumed to be the most important and informative. The goal is to capture as much of this spread as possible in as few dimensions as possible. This is the more common and intuitive explanation of PCA.

### Reconstruction Error Minimization View

- **Concept:** This view focuses on finding a low-dimensional representation that is the **most faithful approximation** of the original data.
- **Objective:** Find a  $k$ -dimensional linear subspace such that the average squared distance between the original data points and their projections onto this subspace is minimized.
- **Intuition:** Imagine you have a 3D cloud of points and want to represent it with a 2D plane. You would position the plane such that it cuts through the "middle" of the cloud, minimizing the average distance from each point to the plane. The axes defining this optimal plane are the principal components.

### Equivalence

It can be proven that these two objectives are two sides of the same coin. A key part of the proof relies on the Pythagorean theorem in high-dimensional space. The total variance of a data point can be decomposed into two parts:

1. The variance of its projection onto the subspace ( $d_1^2$ ).
2. The squared distance from the point to the subspace (the reconstruction error,  $d_2^2$ ).

$$\text{Total Variance} = d_1^2 + d_2^2$$

Since the total variance is fixed, **maximizing the projected variance ( $d_1^2$ ) is mathematically equivalent to minimizing the reconstruction error ( $d_2^2$ )**. Choosing the subspace that does one automatically achieves the other.

---

## Question 25

Derive PCA via Singular Value Decomposition.

### Theory

While PCA is often introduced via eigen-decomposition of the covariance matrix, its modern, numerically stable implementation is based on the Singular Value Decomposition (SVD) of the centered data matrix  $X$ . This derivation shows how the components of SVD directly relate to the principal components and the projected data.

## SVD Definition

Any  $n \times p$  matrix  $X$  (with  $n$  samples,  $p$  features) can be factorized as:

$$X = U\Sigma V^T$$

- $U$ :  $n \times n$  orthogonal matrix of left-singular vectors.
- $\Sigma$ :  $n \times p$  diagonal matrix of singular values  $\sigma_i$ .
- $V$ :  $p \times p$  orthogonal matrix of right-singular vectors.

## Derivation

Let  $X$  be the centered data matrix.

1. **Goal of PCA:** Find the principal components, which are the eigenvectors of the covariance matrix  $C = (1/(n-1))X^T X$ .
2. **Relate Covariance to SVD:** Let's substitute the SVD of  $X$  into the covariance formula:  
$$C = (1/(n-1)) * (U\Sigma V^T)^T (U\Sigma V^T)$$
  
$$C = (1/(n-1)) * (V\Sigma^T U^T)(U\Sigma V^T)$$
  
Since  $U$  is orthogonal,  $U^T U = I$  (the identity matrix).  
$$C = (1/(n-1)) * V(\Sigma^T \Sigma)V^T$$
3. **Identify the Eigen-decomposition:**
  - The matrix  $\Sigma^T \Sigma$  is a  $p \times p$  diagonal matrix whose diagonal entries are the squared singular values  $\sigma_1^2, \sigma_2^2, \dots$ . Let's call this matrix  $S$ .
  - The equation now is  $C = (1/(n-1)) * V S V^T$ . Rearranging gives:  
$$CV = (1/(n-1)) * VS$$
  - This equation has the exact form of an eigen-decomposition  $CV = V\Lambda$ , where  $V$  is the matrix of eigenvectors and  $\Lambda$  is the diagonal matrix of eigenvalues.
  - By comparing the two forms, we can see that:
    - The columns of  $V$  (the right-singular vectors of  $X$ ) are the **eigenvectors** of the covariance matrix  $C$ . Therefore, **the right-singular vectors  $V$  are the principal components**.
    - The diagonal matrix of eigenvalues  $\Lambda$  is equal to  $(1/(n-1))S$ . This means the eigenvalues are  $\lambda_i = \sigma_i^2 / (n-1)$ . The **eigenvalues are proportional to the squared singular values**.
  -
- 4.
5. **Find the Projected Data (Scores):** The projected data  $T$  (the scores) is  $X$  projected onto the principal components  $V$ .  
$$T = XV = (U\Sigma V^T)V$$
  
Since  $V$  is orthogonal,  $V^T V = I$ .  
$$T = U\Sigma$$
  
This shows that the **scores matrix is simply the product of the left-singular vectors and the singular values**.

This derivation provides a complete solution for PCA directly from the SVD of the data matrix, which is more direct and numerically stable.

---

## Question 26

**Discuss importance of data centering before PCA.**

### Theory

Data centering, which involves subtracting the mean of each feature from all of its values, is a **mandatory preprocessing step** for PCA. Failing to center the data will lead to incorrect and uninterpretable results because the principal components will be calculated relative to the origin (zero) instead of the data's center of mass.

### Explanation of Importance

1. **PCA's Objective is About Variance:** PCA's goal is to find the directions of maximum **variance**. Variance, by its statistical definition, is the measure of spread *around the mean*. If the data is not centered, the "mean" that PCA implicitly uses is the origin (0, 0, ..., 0).
2. **Incorrect First Principal Component:**
  - Imagine a cluster of data points located far from the origin. If this data is not centered, the first principal component will not point in the direction of the data's maximum spread. Instead, it will simply point from the origin towards the center of the data cloud.
  - This first component would be capturing the **location** of the data in the feature space, not its **variance structure**, which is what we are actually interested in. It becomes a meaningless component that reflects the arbitrary choice of origin.
- 3.
4. **Covariance Calculation:** The formula for the covariance between two variables A and B is  $\text{Cov}(A, B) = E[(A - \mu_A)(B - \mu_B)]$ . This formula inherently involves subtracting the means. When we compute the covariance matrix  $(1/(n-1))X^T X$ , it is implicitly assumed that X is the centered data matrix where the column means are zero. If X is not centered, this calculation will not yield the true covariance matrix, and the subsequent eigen-decomposition will be incorrect.

### The Effect of Not Centering

- The first principal component will be heavily skewed and will not represent the true direction of maximum variance.
- All subsequent components will also be calculated incorrectly as they must be orthogonal to the first.
- The interpretation of the components becomes impossible.

**Conclusion:** Centering the data ensures that PCA focuses solely on capturing the variance structure of the data relative to its mean. It is a non-negotiable step in the PCA pipeline. Note that standardization (StandardScaler) includes centering by default.

---

## Question 27

### Why does scaling features change component loadings?

#### Theory

Scaling features changes the component loadings because PCA is a variance-based technique. The loadings represent the contribution of each original feature to a principal component. When we scale the data, we fundamentally alter the variance of each feature, which in turn changes their relative importance in the PCA calculation and thus changes the loadings.

#### Explanation

1. **PCA is Variance-Driven:** PCA identifies the first principal component as the direction that explains the maximum possible variance in the data. Subsequent components explain the maximum remaining variance.
2. **Unscaled Data:** In unscaled data, a feature with a large numerical range (e.g., income in dollars) will have a much larger variance than a feature with a small range (e.g., number\_of\_children). PCA will be biased towards the high-variance feature. The first principal component (PC1) will be almost entirely aligned with the income axis. Consequently, the **loading** of income on PC1 will be very high (close to 1), while the loading of number\_of\_children will be very low (close to 0).
3. **Scaled Data (Standardization):** When we standardize the data, we transform every feature to have a variance of 1.
  - This equalizes the influence of all features. No single feature can dominate the analysis due to its original scale.
  - Now, PCA will find the directions of maximum variance based on the **correlations** between the features, not their arbitrary scales.
  - The resulting principal components will be different. PC1 will now be a combination of features that are most strongly correlated. The **loadings will change dramatically**. income and number\_of\_children will now have loadings that reflect their contribution to the shared variance structure, not just their initial scale.
- 4.

#### In summary:

- **Before scaling:** Loadings are influenced by both the **correlation structure** and the **scale** of the original variables.

- **After scaling:** Loadings are influenced only by the **correlation structure**.

Since the goal of PCA is typically to uncover the underlying correlation patterns, scaling is essential for obtaining meaningful and interpretable component loadings.

---

## Question 28

**Explain meaning of scree plot and elbow criteria.**

### Theory

A **scree plot** is a fundamental diagnostic tool used in PCA and Factor Analysis to help determine the optimal number of components or factors to retain. It is a line plot or bar chart that displays the eigenvalues (representing the explained variance) of each component, sorted in descending order.

The **elbow criterion** is a heuristic applied to the scree plot to find a natural cut-off point for the number of components.

### Explanation

#### 1. The Scree Plot:

- **X-axis:** The principal component number (e.g., 1, 2, 3, ...).
- **Y-axis:** The corresponding eigenvalue for each principal component.
- **Interpretation:** The plot shows the magnitude of variance captured by each successive component. Since eigenvalues are sorted in descending order, the plot will always be downward sloping. The steepness of the initial slope shows that the first few components capture a large amount of variance, while the later components capture progressively less.

2.

#### 3. The Elbow (or "Scree") Criterion:

- **Concept:** The name "scree" comes from a geology term for the pile of rubble that accumulates at the bottom of a cliff. In the scree plot, the "cliff" is the steep initial drop, and the "scree" is the flat-ish tail of the plot where the eigenvalues are small and similar to each other.
- **The "Elbow":** This is the point on the plot where the curve bends and begins to flatten out. It marks the point of diminishing returns.
- **Criterion:** The rule suggests that you should retain all the components **before the elbow** (the "cliff" part), as these are the ones that capture a significant and meaningful amount of variance. The components after the elbow (the "scree" part) are considered to represent minor variations or noise and can be discarded.

4.

## Use Case

The primary use of the scree plot and elbow criterion is to provide a visual, data-driven heuristic for choosing the number of components ( $k$ ) to keep for dimensionality reduction. It helps answer the question, "How many components are needed to represent the important signal in the data without including too much noise?"

## Pitfalls

- The elbow is often not sharp and can be ambiguous, making the choice subjective.
  - It should be used in conjunction with other methods, like the cumulative explained variance ratio, to make a more robust decision.
- 

## Question 29

**How to decide number of components via Kaiser rule.**

### Theory

The **Kaiser Rule**, also known as the Kaiser-Guttman criterion, is another popular heuristic for determining the number of principal components or factors to retain. It is a simple, objective rule that is easy to apply.

### The Rule Explained

The Kaiser Rule states that you should **retain only those principal components whose eigenvalues are greater than 1.0**.

### Rationale and Context

1. **Standardized Data:** This rule should only be applied when PCA is performed on **standardized data**. When data is standardized, each original feature has a variance of 1.0. The total variance in the dataset is therefore equal to the number of features ( $p$ ).
2. **Average Variance:** The average eigenvalue for a dataset with  $p$  features will be  $(\text{Total Variance}) / p = p / p = 1.0$ .
3. **The Logic:** The Kaiser Rule's logic is that any principal component with an eigenvalue less than 1.0 is explaining less variance than a single original variable did. Therefore, such a component is not considered significant or useful, as it is less informative than an individual standardized feature. Retaining it would not be an efficient way to summarize the data.

### How to Apply

1. Standardize the data.

2. Perform PCA and obtain the eigenvalues for all components.
3. Count how many of these eigenvalues are greater than 1.0.
4. This count is the number of components to retain according to the rule.

### Pitfalls and Criticism

- **Arbitrariness:** The 1.0 cut-off is arbitrary and has been criticized for being too strict or too lenient depending on the dataset.
  - **Over- or Underestimation:** For datasets with a small number of variables, it can sometimes suggest keeping too few components. For datasets with a very large number of variables, it can suggest keeping too many.
  - **Not a Replacement for Judgment:** Like the elbow criterion, the Kaiser Rule is a heuristic, not a definitive law. It is best used as a starting point or as one piece of evidence alongside other methods like the scree plot and cumulative explained variance. Many researchers prefer the scree plot because it is based on the observed data structure rather than a fixed, arbitrary value.
- 

## Question 30

**Describe cumulative explained variance ratio.**

### Theory

The **cumulative explained variance ratio** is a diagnostic measure used in PCA to determine how much of the total information (variance) in the original dataset is captured by a given number of principal components. It is one of the most common and effective methods for deciding how many components to retain for dimensionality reduction.

### Explanation

1. **Explained Variance Ratio:**
  - First, we calculate the **explained variance ratio** for each individual principal component  $i$ . This is the proportion of the total variance that this single component accounts for.
  - $\text{Ratio}(\text{PC}_i) = \text{Eigenvalue}(\text{PC}_i) / \text{Sum of all Eigenvalues}$
- 2.
3. **Cumulative Calculation:**
  - The cumulative explained variance ratio for the first  $k$  components is the sum of their individual explained variance ratios.
  - $\text{Cumulative Ratio}(k) = \text{Ratio}(\text{PC}_1) + \text{Ratio}(\text{PC}_2) + \dots + \text{Ratio}(\text{PC}_k)$
- 4.

### How It's Used

1. **Set a Threshold:** The primary use is to decide on the number of components ( $k$ ) to keep. A data scientist will typically set a threshold for the desired amount of variance to preserve, for example, 90%, 95%, or 99%.
2. **Find the Number of Components:** They then calculate the cumulative explained variance for an increasing number of components until the chosen threshold is met or exceeded.
  - o **Example:** If the cumulative variance for the first 5 components is 0.88 (88%) and for the first 6 components is 0.91 (91%), and the goal was to retain 90% of the variance, then  $k=6$  components would be selected.
- 3.

## Visualization

This is often visualized by overlaying a cumulative explained variance curve on top of a scree plot.

- The scree plot shows the variance of individual components.
- The cumulative curve shows the total variance captured as we add more components.
- The user can then see, for example, that they need to keep 10 components to cross the 95% threshold.

## Best Practices

- The choice of threshold (e.g., 95%) is domain-specific. For some applications where data fidelity is critical, a 99% threshold might be needed. For others where only the main patterns matter, 80% might suffice.
  - This method provides a more objective and interpretable way to choose  $k$  compared to the subjective elbow criterion.
- 

## Question 31

**Explain whitening transformation and its risks.**

### Theory

**Whitening**, in the context of PCA, is a transformation applied to the data *after* the standard PCA projection. Its goal is to produce a new set of variables that are not only **uncorrelated** (which PCA already achieves) but also have a **unit variance**. This means the resulting covariance matrix of the whitened data is the identity matrix.

### How Whitening Works

1. **Perform Standard PCA:** First, center the data and compute the principal component scores  $T = XV$ , where  $X$  is the centered data and  $V$  are the principal components (eigenvectors).
2. **Scale by Eigenvalues:** The variance of each principal component score is its corresponding eigenvalue  $\lambda_i$ . To give each component a unit variance, we simply divide each component score by the square root of its eigenvalue.
  - o Whitened Score  $T_{\text{whitened}}(i) = T(i) / \sqrt{\lambda_i}$
- 3.

The resulting  $T_{\text{whitened}}$  matrix contains the whitened data.

### **Benefits of Whitening**

- **Equal Importance:** Whitening makes all the transformed features have the same variance. This can be beneficial for some machine learning algorithms (like some types of clustering or neural networks) that assume the input features are on a similar scale and are uncorrelated.
- **Preprocessing Step:** It can be seen as a more thorough form of preprocessing, creating a simplified, spherical data structure.

### **Risks and Disadvantages**

1. **Loss of Information:** This is the primary risk. The explained variance of each principal component is a crucial piece of information. The first component is "principal" precisely because it has the highest variance. Whitening destroys this information by forcing all components to have the same variance (1.0). You no longer know which components were originally the most important.
2. **Noise Amplification:** PCA is often used for noise reduction by discarding low-variance components. Whitening does the opposite for these components. It takes the low-variance components (which are often noise) and scales them *up* so that their variance becomes 1.0. This can significantly amplify the noise in the data, which may be detrimental to the performance of a subsequent model.

**Conclusion:** Whitening should be used with extreme caution. It is useful in a few specific contexts (e.g., as a step in certain statistical algorithms like Independent Component Analysis (ICA)), but for general-purpose dimensionality reduction or feature engineering, it is usually not recommended because it discards the valuable information contained in the relative variances of the principal components.

---

## **Question 32**

**Compare PCA with Factor Analysis.**

## Theory

PCA and Factor Analysis (FA) are both linear dimensionality reduction techniques that are often confused, but they have fundamentally different assumptions, goals, and interpretations. PCA is a mathematical data-simplification technique, while FA is a statistical model for uncovering latent structure.

### Comparison Table

Feature	Principal Component Analysis (PCA)	Factor Analysis (FA)
<b>Primary Goal</b>	To account for the <b>total variance</b> in the observed variables. It finds components that best summarize the data.	To explain the <b>common variance</b> (covariance) among observed variables by identifying unobserved <b>latent factors</b> .
<b>Underlying Model</b>	<b>No formal model.</b> PCA is an algorithm that transforms the data. The principal components are linear combinations of the observed variables.	<b>Assumes a formal latent variable model:</b> Observed_Variable = Loading * Factor + Error. It assumes the correlations between observed variables are due to these hidden factors.
<b>Variance Decomp.</b>	It decomposes the total variance of the variables.	It partitions the variance of each variable into two parts: <b>common variance</b> (explained by the latent factors) and <b>unique variance</b> (error and variable-specific variance).
<b>Relationship</b>	The components are defined as a function of the variables: PC = f(Variables).	The variables are defined as a function of the factors: Variable = f(Factors, Error).
<b>Output</b>	Principal components are orthogonal. The key outputs are the component scores and the explained variance.	Factors are estimated and are not necessarily orthogonal (can be rotated). The key outputs are <b>factor loadings</b> , which are the correlations between variables and factors.

### When to Use Which?

- **Use PCA when your goal is purely data reduction.** You want to reduce the number of features in a way that preserves as much information as possible, for example, to speed up a machine learning model or for visualization.
- **Use Factor Analysis when your goal is interpretation and theory.** You want to understand the underlying constructs or latent variables that are causing your observed

variables to be correlated. This is common in social sciences, psychology, and marketing to develop and validate theories.

---

## Question 33

**Describe kernel PCA and nonlinear embeddings.**

### Theory

**Kernel PCA (kPCA)** is a powerful extension of PCA designed to handle **non-linear** relationships in data. It achieves this by using the "kernel trick" to implicitly map the data into a higher-dimensional feature space where the non-linear patterns may become linear, and then performs PCA in that space.

A **non-linear embedding** is a low-dimensional representation of data that preserves its underlying non-linear manifold structure. Kernel PCA is one method for creating such an embedding.

### How Kernel PCA Works

1. **The Limitation of Linear PCA:** Standard PCA can only find linear subspaces. It fails on data with non-linear structures like spirals or concentric circles.
2. **The Kernel Trick:** The key idea is to avoid the explicit, computationally expensive mapping of the data into a high-dimensional space  $\Phi(x)$ . Instead, all the operations of PCA can be rewritten in terms of the dot products of the data points,  $x_i^T x_j$ . The kernel trick replaces this dot product with a **kernel function**,  $K(x_i, x_j) = \Phi(x_i)^T \Phi(x_j)$ . This calculates the dot product in the high-dimensional space without ever having to go there.
3. **Common Kernels:**
  - **Polynomial Kernel:** Captures polynomial relationships.
  - **Radial Basis Function (RBF) Kernel:** A very popular choice that can handle complex, non-linear structures.
  - **Sigmoid Kernel:** Similar to the activation function in neural networks.
- 4.
5. **The Process:** kPCA computes the  $n \times n$  kernel matrix  $K$  for all pairs of data points. It then performs an eigen-decomposition on this kernel matrix to find the principal components in the high-dimensional feature space. The data is then projected onto these components to create the non-linear embedding.

### Advantages and Disadvantages

- **Advantage:** Can effectively capture and reduce the dimensionality of non-linear data structures.

- **Disadvantage:** It is computationally more expensive ( $O(n^2)$ ) than standard PCA, is sensitive to the choice of kernel and its parameters (e.g., gamma for RBF), and the resulting components are not easily interpretable.
- 

## Question 34

Discuss incremental PCA for streaming data.

### Theory

**Incremental PCA (IPCA)** is a variant of PCA specifically designed for scenarios where the dataset is too large to fit in memory or when the data arrives in a continuous stream. It processes the data in mini-batches, allowing it to learn the principal components in an online or incremental fashion without needing the entire dataset at once.

### Application to Streaming Data

In a streaming data context, data points arrive one by one or in small chunks. Retraining a standard PCA model from scratch every time new data arrives would be computationally infeasible. IPCA provides an efficient solution.

1. **Model Initialization:** An IPCA model is initialized, either with no prior information or by training on an initial chunk of data.
2. **Processing a New Batch:** As a new batch of streaming data arrives:
  - The `partial_fit()` method of the IPCA model is called with the new batch.
  - The model updates its current estimates of the principal components (eigenvectors) and explained variances (eigenvalues) based on this new data and the summary statistics it has maintained from previous batches.
- 3.
4. **Continuous Updates:** This process is repeated for each incoming batch. The model continuously adapts and refines its understanding of the data's variance structure as more data is observed.

### Trade-offs and Considerations

- **Approximation:** The result of IPCA is an approximation of the true principal components that would be found by running batch PCA on the entire stream up to that point. The quality of the approximation depends on the batch size.
- **Concept Drift:** Standard IPCA assumes the underlying data distribution is stationary. If the statistical properties of the stream change over time (**concept drift**), the model may become outdated. More advanced online PCA methods exist that can "forget" older data to adapt to such drifts.

- **Memory Efficiency:** Its key advantage is its low memory footprint, which only depends on the size of the mini-batch, not the total size of the stream.

**Use Case:** Monitoring network traffic data in real-time. IPCA could be used to continuously find the main patterns of traffic variation, helping to detect anomalies or changes in network behavior as they happen.

---

## Question 35

**Explain randomized SVD acceleration.**

### Theory

**Randomized SVD** is a computational technique used to speed up the Singular Value Decomposition (SVD) of a matrix, particularly for large matrices. Since PCA is often implemented using SVD, randomized SVD can be used to create a much faster, approximate version of PCA. It is the core of scikit-learn's PCA(`svd_solver='randomized'`) option.

### The Core Idea

The main idea behind randomized SVD is that for a large matrix, its essential structure can be captured by a lower-rank approximation. Instead of computing the full, exact SVD, the algorithm finds this approximation using a probabilistic approach.

### The Process

1. **Find an Approximate Subspace:** The key step is to find a smaller matrix Q whose range (the space spanned by its columns) approximates the range of the original matrix X. This is done efficiently using random projections.
  - A random matrix R is created.
  - The product Y = XR is computed.
  - An orthonormal basis Q is found for Y (e.g., using QR decomposition). This Q now captures most of the action of X.
- 2.
3. **Project and Compute SVD:** The original large matrix X is projected onto this smaller subspace:  $B = Q^T X$ . This projected matrix B is much smaller than X.
4. **Compute SVD on the Small Matrix:** A standard, exact SVD is now computed on the small matrix B.
5. **Reconstruct the Final SVD:** The SVD of the original matrix X is then approximated from the SVD of B and the projection matrix Q.

### Advantages

- **Speed:** Randomized SVD is significantly faster than exact SVD for large matrices, especially when the desired number of components  $k$  is much smaller than the number of features. Its complexity is much lower.
- **Accuracy:** It provides a very good approximation of the top singular values and vectors. The results are often very close to what exact PCA would produce.

### When to Use

- This is the ideal choice for **high-dimensional datasets** (where the number of features is large) when you only need to compute the first few principal components.
  - scikit-learn's PCA uses this method by default when  $n\_components$  is small compared to the data size, as it offers the best performance trade-off in these common scenarios.
- 

## Question 36

**Describe robust PCA and handling outliers.**

### Theory

**Robust PCA** is a modification of standard PCA designed to handle datasets that are contaminated with **gross errors or outliers**. Standard PCA is highly sensitive to outliers because its objective function is based on minimizing the sum of squared errors. A single large outlier can significantly skew the principal components.

Robust PCA is based on the assumption that the data matrix  $D$  can be decomposed into two separate components:

$$D = L + S$$

- $L$ : A **low-rank matrix** that represents the true, underlying clean data structure (the principal components).
- $S$ : A **sparse matrix** that represents the outliers or gross errors. "Sparse" means that it has very few non-zero entries.

### How it Works

The goal of Robust PCA is to find the  $L$  and  $S$  matrices that best reconstruct the original data  $D$ . This is framed as a convex optimization problem, often called **Principal Component Pursuit**, which aims to find the  $L$  with the lowest possible rank and the  $S$  with the highest possible sparsity (lowest number of non-zero elements).

This decomposition effectively separates the underlying structure from the sparse outliers.

### Comparison with Standard PCA

- **Sensitivity:** Standard PCA is sensitive to outliers. Robust PCA is, by design, robust to them.
- **Output:** Standard PCA outputs principal components and scores. Robust PCA outputs two distinct matrices: the cleaned low-rank data L and the identified outlier matrix S.
- **Use Case:** Standard PCA is for general dimensionality reduction and feature extraction. Robust PCA is specifically for **recovering the principal components of a dataset in the presence of significant outliers**.

## Applications

- **Video Surveillance:** Used for background subtraction. A video sequence can be represented as a matrix D. The static background is the low-rank component L, and moving objects (like people or cars) are the sparse component S.
  - **Image Denoising:** Can be used to remove "salt and pepper" noise or corrupted pixels from an image.
  - **Data Cleaning:** Can be used as a preprocessing step to identify and separate outliers before further analysis.
- 

## Question 37

**Explain projecting new samples into PCA space.**

### Theory

After a PCA model has been trained on a dataset, it can be used to project new, unseen data samples into the same principal component space. This is a crucial step for ensuring consistency, for example, when applying a machine learning model that was trained on PCA-transformed data to new prediction requests.

### The Process

1. **Train the PCA Model:** First, fit a PCA model on the **training dataset**. This involves:
  - Learning the mean and standard deviation of the training data (for scaling).
  - Learning the principal components (eigenvectors) from the training data.
  - The fit() method in scikit-learn performs these steps.
- 2.
3. **Preprocess the New Data:**
  - Take the new data sample(s).
  - Crucially, you must preprocess the new data using the **exact same parameters** learned from the training data.
  - This means centering the new data by subtracting the **training data's mean** and scaling it by the **training data's standard deviation**. You should not re-calculate the mean or std dev from the new data.

4.

5. **Project the New Data:**

- The final step is to project the preprocessed new data onto the principal components that were learned from the training data.
- This is a matrix multiplication:  $\text{New\_Scores} = \text{New\_Data\_Processed} * V$ , where  $V$  is the matrix of principal components (eigenvectors) learned from the training set.
- In scikit-learn, this is done simply by calling the `transform()` method of the already-fitted PCA object on the new data.

6.

**Code Example (Conceptual)**

```
code Python
downloadcontent_copyexpand_less
IGNORE_WHEN COPYING_START
IGNORE_WHEN COPYING_END
    from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler

# 1. Fit on training data
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)

pca = PCA(n_components=10)
X_train_pca = pca.fit_transform(X_train_scaled)

# --- Now, handle a new data point X_new ---

# 2. Preprocess the new data using the FITTED scaler and pca
X_new_scaled = scaler.transform(X_new) # Use .transform(), NOT .fit_transform()
X_new_pca = pca.transform(X_new_scaled) # Use .transform(), NOT .fit_transform()

# X_new_pca is now the representation of the new data point in the
# same PCA space as the training data.
```

**Pitfalls**

The most common mistake is to fit the scaler or the PCA model on the new data separately. This would create a completely different transformation with a different mean, scale, and set of components, making the new projected data incompatible with any model trained on the original projected data. **Always fit on the training set and transform on the test/new set.**

---

# Question 38

Discuss PCA for missing-value imputation.

## Theory

PCA can be used as a method for imputing missing values, often providing more accurate results than simple mean or median imputation. The approach leverages the correlation structure of the data, as captured by the principal components, to predict the likely values of the missing entries. This is typically implemented using an iterative algorithm.

### Iterative PCA Imputation (EM-like approach)

1. **Initial Imputation:** Start by making an initial, simple guess for the missing values. A common choice is to impute them with the mean of their respective columns.
2. **Iterative Process:** Repeat the following steps until the imputed values converge:
  - **Step A (PCA):** Perform PCA on the current, complete version of the dataset. Select a number of components  $k$  to represent the main data structure.
  - **Step B (Reconstruction):** Transform the data into the  $k$ -dimensional PCA space and then immediately transform it back into the original feature space (`inverse_transform`). This produces a reconstructed, "denoised" version of the data.
  - **Step C (Update):** Replace the values in the original matrix that were initially missing with their corresponding values from the reconstructed matrix. The other values (the ones that were originally present) are left untouched.
- 3.
4. **Convergence:** The algorithm stops when the difference between the imputed values from one iteration to the next falls below a certain threshold.

## Why It Works

This method works because the PCA reconstruction provides a better estimate for the missing values than a simple mean. The reconstruction is based on the dominant correlation patterns in the data. By projecting onto the principal components, the algorithm effectively uses the information from the features that are present for a given sample to predict what the missing feature's value should be to fit the overall data structure.

## Libraries

This iterative approach is implemented in libraries like `sklearn.impute.IterativeImputer` (which can use various estimators, not just PCA) and in specific PCA packages in R (e.g., the `missMDA` package).

## Limitations

- It is computationally more expensive than simple imputation.

- It assumes the data has a low-rank structure that can be well-approximated by a few principal components.
- 

## Question 39

**Compare PCA vs. autoencoders for dimensionality reduction.**

### Theory

PCA and autoencoders are both powerful techniques for dimensionality reduction, but they differ fundamentally in their linearity and how they learn the data representation. PCA is a linear algebraic method, while autoencoders are a type of neural network capable of learning non-linear transformations.

### Comparison Table

Feature	Principal Component Analysis (PCA)	Autoencoder
<b>Linearity</b>	<b>Linear.</b> PCA can only learn a linear projection of the data.	<b>Non-linear.</b> By using non-linear activation functions in its hidden layers, an autoencoder can learn complex, non-linear mappings.
<b>Learning Method</b>	<b>Algorithmic.</b> It is a deterministic mathematical procedure (eigen-decomposition or SVD). No iterative training is needed.	<b>Learned via Training.</b> It is a neural network that must be trained iteratively using backpropagation to minimize reconstruction error.
<b>Parameters</b>	<b>No hyperparameters</b> to tune (other than the number of components to keep).	<b>Many hyperparameters</b> to tune: number of hidden layers, number of neurons per layer, activation functions, learning rate, optimizer, etc.
<b>Performance</b>	Very fast and computationally efficient.	Can be very slow to train, especially for large datasets and complex architectures.
<b>Interpretability</b>	The components are somewhat interpretable through loadings.	The learned encoding is a "black box" and is generally not interpretable.

<b>Data Requirement</b>	Works well even on smaller datasets.	Requires a large amount of data to train effectively and avoid overfitting.
-------------------------	--------------------------------------	---

### When to Use Which?

- **Use PCA as a first-line approach.** It is fast, simple, and provides a great baseline. If the underlying structure of your data is mostly linear, PCA will perform very well.
- **Use an autoencoder when you need more power and suspect non-linearities.** If PCA provides a poor representation (e.g., a supervised model performs badly on the PCA-transformed data), it suggests that non-linear relationships are important. An autoencoder can capture these complex patterns, potentially leading to a much better low-dimensional representation, at the cost of higher complexity and slower training.

**Relationship:** A simple linear autoencoder (an autoencoder with no non-linear activation functions) with a single hidden layer can be shown to be equivalent to PCA. It will learn a latent space that spans the same subspace as the principal components found by PCA.

---

## Question 40

Explain biplot interpretation.

### Theory

A **biplot** is a specialized and highly informative visualization that combines the results of PCA into a single graph. It simultaneously displays information about both the **samples** (the data points) and the **features** (the original variables). The "bi" in biplot refers to this dual representation.

### What a Biplot Shows

A biplot is typically a 2D scatter plot where:

1. **The Samples:** The data points are plotted as points, just like in a standard PCA scatter plot. Their positions are determined by their scores on the first two principal components (PC1 and PC2).
2. **The Features:** The original features are plotted as **vectors** (arrows) originating from the center of the plot. These vectors represent the loadings of the features on the principal components.

### How to Interpret a Biplot

1. **Interpreting the Samples (Points):**
  - Points that are close to each other are similar in their overall feature profiles.

- Clusters of points indicate potential groups or classes in the data.
- 2.
- 3. **Interpreting the Features (Vectors):**
  - **Vector Length:** The length of a feature's vector indicates its magnitude or influence on the plotted space. Longer vectors represent features that are better explained by the first two components.
  - **Vector Direction:** The direction of the vector shows how the feature contributes to the principal components. A vector pointing mostly along the PC1 axis means that feature has a high loading on PC1.
- 4.
- 5. **Interpreting the Relationship Between Features:**
  - **Angle Between Vectors:** The cosine of the angle between two feature vectors is proportional to the **correlation** between those two original features.
    - **Small Angle ( $\sim 0^\circ$ ):** The features are highly positively correlated.
    - **Large Angle ( $\sim 180^\circ$ ):** The features are highly negatively correlated.
    - **Right Angle ( $\sim 90^\circ$ ):** The features are uncorrelated.
  -
- 6.
- 7. **Interpreting the Relationship Between Samples and Features:**
  - Projecting a sample point onto a feature's vector gives an approximation of that sample's value for that feature. If a sample point lies far along the direction of a feature vector, it means that sample has a high value for that particular feature.
- 8.

A biplot provides a rich, holistic view of the data, allowing a data scientist to understand the structure of the samples, the importance of the features, and the relationships between them all in a single plot.

Here are the completed answers for all the questions.

## Pca Interview Questions - General Questions

### Question 1

**How is PCA used for dimensionality reduction?**

#### Theory

PCA is used for dimensionality reduction by transforming a large set of correlated variables into a smaller set of uncorrelated variables, called principal components, while preserving the

maximum amount of information from the original data. This "information" is quantified as the variance in the dataset.

## Explanation

The process involves creating a new coordinate system for the data, where the axes of this new system (the principal components) are aligned with the directions of maximum variance.

1. **Find Principal Components:** PCA identifies the directions (principal components) in the data that capture the most variance. The first principal component (PC1) is the direction that explains the largest possible variance. The second principal component (PC2) is orthogonal (perpendicular) to PC1 and explains the largest remaining variance, and so on.
2. **Select a Subset of Components:** Each principal component has a corresponding eigenvalue, which represents the amount of variance it captures. By sorting these eigenvalues in descending order, we can rank the components by their importance. To reduce dimensionality, we select only the top  $k$  components that capture a significant portion of the total variance (e.g., enough components to explain 95% of the total variance).
3. **Project the Data:** The original high-dimensional data is then projected onto the lower-dimensional subspace defined by these selected  $k$  principal components. The result is a new dataset with  $k$  features (where  $k$  is much smaller than the original number of features), which is a lower-dimensional representation of the original data.

## Use Cases

- **Speeding up Machine Learning Models:** Reducing the number of features significantly decreases the computational cost of training a model.
  - **Combating the Curse of Dimensionality:** By reducing dimensions, PCA helps mitigate issues like data sparsity and overfitting that arise in high-dimensional spaces.
  - **Data Visualization:** Projecting data onto 2 or 3 principal components allows for visualization of high-dimensional data structures in 2D or 3D scatter plots.
- 

## Question 2

**Why is PCA considered an unsupervised technique?**

### Theory

PCA is considered an **unsupervised** technique because it operates on the input data  $X$  alone, without any knowledge of, or reference to, a target variable or labels ( $y$ ). Its goal is to find the intrinsic structure (the directions of maximum variance) within the data itself, not to predict an output.

## Explanation

1. **No Labels Required:** The entire PCA algorithm—standardizing the data, computing the covariance matrix, and finding eigenvectors/eigenvalues—only uses the feature values of the dataset. It does not look at any labels that might be associated with the data points.
2. **Objective is Data-Internal:** The objective of PCA is to find a lower-dimensional representation that maximizes variance or minimizes reconstruction error. This objective is defined solely by the relationships between the features within the dataset. It is not trying to find features that are predictive of a specific outcome.
3. **Contrast with Supervised Techniques:**
  - **Supervised dimensionality reduction techniques**, like Linear Discriminant Analysis (LDA), explicitly use the class labels to find a projection that maximizes the separability between classes. LDA's objective is supervised.
  - **Supervised feature selection methods**, like Recursive Feature Elimination, evaluate features based on how well they contribute to the performance of a model in predicting a target variable.
- 4.

## Use Cases

Because it is unsupervised, PCA is a general-purpose tool that can be used as a preprocessing step for both supervised learning (classification, regression) and unsupervised learning (clustering) tasks. It finds the inherent structure of the features before they are used for any downstream task.

---

## Question 3

**Derive the PCA from the optimization perspective, i.e., minimization of reconstruction error.**

### Theory

PCA can be derived from the objective of finding a low-dimensional projection of the data that minimizes the reconstruction error. This means we want to find a k-dimensional subspace such that when the data is projected onto it and then projected back into the original space, the average squared distance between the original points and their reconstructions is as small as possible.

### Mathematical Derivation

Let  $X$  be the  $n \times p$  centered data matrix. We want to project the data onto a k-dimensional subspace defined by an orthonormal basis  $V$  ( $a p \times k$  matrix where  $V^T V = I$ ).

1. **Projection and Reconstruction:**
  - A data point  $x_i$  is projected into the subspace as  $V\mathbb{Q}^T x_i$ .
  - The reconstruction  $\hat{x}_i$  is obtained by projecting it back:  $\hat{x}_i = V\mathbb{Q}(V\mathbb{Q}^T x_i)$ .
- 2.
3. **Objective Function:** We want to minimize the sum of squared reconstruction errors  $J$ :  

$$J = \sum_i \|x_i - \hat{x}_i\|^2 = \sum_i \|x_i - V\mathbb{Q}V\mathbb{Q}^T x_i\|^2$$
4. **Simplifying the Objective:** Using the property  $\|a - b\|^2 = a^T a - 2a^T b + b^T b$  and properties of orthonormal matrices:  

$$\|x_i - V\mathbb{Q}V\mathbb{Q}^T x_i\|^2 = x_i^T x_i - 2x_i^T V\mathbb{Q}V\mathbb{Q}^T x_i + x_i^T V\mathbb{Q}V\mathbb{Q}^T V\mathbb{Q}V\mathbb{Q}^T x_i$$
  
Since  $V\mathbb{Q}^T V\mathbb{Q} = I$ , this simplifies to:  

$$= x_i^T x_i - x_i^T V\mathbb{Q}V\mathbb{Q}^T x_i$$
5. **Connecting to Variance:** The term  $x_i^T x_i$  is related to the total variance and is constant with respect to  $V\mathbb{Q}$ . So, minimizing  $J$  is equivalent to **maximizing** the term  $\sum_i x_i^T V\mathbb{Q}V\mathbb{Q}^T x_i$ . This can be rewritten using the trace operator:  

$$\sum_i x_i^T V\mathbb{Q}V\mathbb{Q}^T x_i = \text{Tr}(V\mathbb{Q}^T (\sum_i x_i x_i^T) V\mathbb{Q}) = \text{Tr}(V\mathbb{Q}^T (X^T X) V\mathbb{Q})$$
6. **Final Optimization Problem:** We need to maximize  $\text{Tr}(V\mathbb{Q}^T C V\mathbb{Q})$ , where  $C$  is proportional to the covariance matrix  $X^T X$ , subject to  $V\mathbb{Q}$  being orthonormal.
  - It is a standard result in linear algebra that this trace is maximized when the columns of  $V\mathbb{Q}$  are the  $k$  eigenvectors of the matrix  $C$  corresponding to the  $k$  largest eigenvalues.
- 7.

**Conclusion:** The subspace that minimizes the reconstruction error is the same one spanned by the top  $k$  principal components (the eigenvectors of the covariance matrix), which is the same solution obtained from the variance maximization perspective.

---

## Question 4

**How do you determine the number of principal components to use?**

### Theory

Choosing the number of principal components ( $k$ ) to retain is a critical decision in PCA that involves a trade-off between dimensionality reduction and information preservation. There is no single "best" method; a combination of techniques is often used to make a well-informed choice.

### Multiple Solution Approaches

1. **Cumulative Explained Variance Ratio:**
  - **Method:** This is the most common and interpretable method. Calculate the cumulative sum of the explained variance ratios of the components (sorted by importance).

- **Procedure:** Choose the smallest number of components  $k$  that captures a desired percentage of the total variance, typically between 90% and 99%. For example, select  $k$  such that the cumulative explained variance is  $\geq 0.95$ .
    - **Best for:** When you have a clear goal for information retention.
- 2.
3. **Scree Plot and Elbow Criterion:**
    - **Method:** Plot the eigenvalues (explained variance) for each component in descending order.
    - **Procedure:** Look for the "elbow" in the plot—the point where the curve flattens out. The number of components to keep is the number at this elbow point. The idea is that components after the elbow contribute minimally and are likely noise.
    - **Best for:** Visual exploration and getting a quick heuristic sense of the data's intrinsic dimensionality.
- 4.
5. **Kaiser Rule:**
    - **Method:** A simple heuristic applied to PCA on **standardized data**.
    - **Procedure:** Retain only those components whose eigenvalues are greater than 1.0. The rationale is that a component should explain at least as much variance as a single original standardized variable.
    - **Best for:** A quick, objective first guess, but it's often criticized for being arbitrary.
- 6.
7. **Application-Specific Performance:**
    - **Method:** If PCA is being used as a preprocessing step for a supervised learning model, you can treat the number of components  $k$  as a hyperparameter.
    - **Procedure:** Use cross-validation to evaluate the performance of the downstream model for different values of  $k$ . Choose the  $k$  that results in the best model performance (e.g., highest accuracy or lowest error).
    - **Best for:** When the ultimate goal is to optimize the performance of a predictive model.
- 8.

**Best Practice:** Combine methods. Start with a scree plot and the cumulative explained variance to identify a reasonable range for  $k$ . If the end goal is a predictive model, fine-tune  $k$  using cross-validation within that range.

---

## Question 5

**Provide examples of how PCA can be used in image processing.**

**Theory**

In image processing, PCA is a versatile tool used for dimensionality reduction, feature extraction, and data compression. An image can be treated as a high-dimensional data vector, making it a prime candidate for PCA.

## Examples of Applications

### 1. Facial Recognition (Eigenfaces):

- **Concept:** This is the classic application of PCA in image processing. A large dataset of facial images is used. Each image is flattened into a long vector of pixel values.
- **Application:** PCA is applied to this dataset. The resulting principal components are called "**eigenfaces**." Each eigenface represents a fundamental "feature" or pattern of variation in human faces (e.g., variation in lighting, nose shape, etc.).
- **Use:** Any face can then be represented as a linear combination of these eigenfaces. To recognize a new face, it is projected into this "face space," and the closest known face in that space is identified.

2.

### 3. Image Compression:

- **Concept:** An image contains a lot of redundant information. PCA can find a more efficient representation.
- **Application:** The image can be divided into smaller blocks (e.g., 8x8 pixels). PCA can be applied to a large collection of these blocks to find the principal components that best represent them. To compress an image, each block is represented by its scores on the top few components.
- **Use:** By storing only the component scores and the principal components themselves, the image can be reconstructed with minimal loss of quality while using significantly less storage space.

4.

### 5. Hyperspectral Image Analysis:

- **Concept:** Hyperspectral images capture data over hundreds of narrow spectral bands (dimensions). These bands are often highly correlated.
- **Application:** PCA is used to reduce the dimensionality of the spectral data from hundreds of bands down to a few principal components. This makes the data easier to process and classify (e.g., for land use classification or crop health monitoring) while retaining the most important spectral information.

6.

### 7. Handwriting Recognition:

- **Concept:** Similar to eigenfaces, PCA can be applied to a dataset of handwritten characters.
- **Application:** It extracts the principal components of variation in writing styles for each character. A new handwritten character can be classified by projecting it onto these components and finding the best match.

8.

---

# Question 6

Can PCA be applied to categorical data? Why or why not?

## Theory

Standard PCA **cannot be directly applied to categorical data**. The reason is fundamental to how PCA works: it is a mathematical technique based on variance and covariance, which are concepts defined for continuous numerical variables.

## Why Standard PCA Fails on Categorical Data

1. **No Meaningful Mean or Variance:** PCA's first step is to center the data by subtracting the mean. For categorical data (e.g., "red," "green," "blue"), the concepts of mean, variance, and covariance are not mathematically defined.
2. **Incorrect Distance Calculation:** Even if categories are numerically encoded (e.g., red=1, green=2, blue=3), this imposes an arbitrary ordinal relationship and scale. PCA would treat the distance between "red" and "blue" as twice the distance between "red" and "green," which is meaningless.
3. **Linearity Assumption:** PCA finds linear combinations of features. Creating linear combinations of numerically encoded categories is not a valid or interpretable operation.

## Solutions for Handling Categorical Data

If you need to perform dimensionality reduction on a dataset containing categorical features, you must use alternative methods:

1. **One-Hot Encoding + PCA:**
  - **Method:** Convert the categorical features into a set of binary (0/1) features using one-hot encoding. You can then apply standard PCA to the resulting numerical matrix.
  - **Limitation:** This can lead to a very high-dimensional and sparse matrix, and the interpretation of the components can still be difficult.
- 2.
3. **Multiple Correspondence Analysis (MCA):**
  - **Concept:** MCA is specifically designed for categorical data. It can be thought of as the **equivalent of PCA for categorical variables**.
  - **Method:** It analyzes the relationships within a contingency table (a table of frequencies) to find a low-dimensional representation of the data that best captures the associations between categories.
- 4.
5. **Factor Analysis of Mixed Data (FAMD):**
  - **Concept:** This technique is designed for datasets that contain a mix of both **numerical and categorical** features.

- **Method:** It handles the two types of variables simultaneously, applying PCA-like logic to the numerical features and MCA-like logic to the categorical features within a unified framework.

6.

**Conclusion:** Do not use standard PCA on raw or label-encoded categorical data. Use MCA for purely categorical data or FAMD for mixed data types.

---

## Question 7

**How can PCA be parallelized to handle very large datasets?**

### Theory

Standard PCA, when implemented via a full SVD, can be computationally intensive for very large datasets. Parallelizing PCA is crucial for handling Big Data. The strategies for parallelization depend on the bottleneck, which is typically the computation of the covariance matrix or the SVD itself.

### Parallelization Strategies

#### 1. Parallel Covariance Matrix Calculation:

- **Concept:** If the number of features  $p$  is small enough that the  $p \times p$  covariance matrix fits in memory, but the number of samples  $n$  is massive, the calculation of the covariance matrix  $X^T X$  can be parallelized.
- **Method:** This is a classic MapReduce-style problem.
  - **Map Step:** The dataset  $X$  is partitioned and distributed across multiple worker nodes. Each worker computes the outer product  $x_i x_i^T$  for its assigned data points  $x_i$ .
  - **Reduce Step:** The results from all workers are summed together to form the final  $X^T X$  matrix.
- Once the distributed covariance matrix is formed, a standard eigen-decomposition can be performed on a single machine (if it fits in memory).

2.

#### 3. Parallel Singular Value Decomposition (SVD):

- **Concept:** For datasets where both  $n$  and  $p$  are large, the SVD itself needs to be parallelized.
- **Method:** This is more complex and involves distributed block-based matrix algorithms.
  - The data matrix  $X$  is partitioned into blocks and distributed across a grid of processors.

- Algorithms like **Tall Skinny QR (TSQR)** followed by a direct SVD on the small R matrix, or iterative methods like the Lanczos algorithm, are used to compute the SVD in a distributed manner.
  - Libraries like ScaLAPACK and frameworks like Apache Spark's MLlib have built-in implementations of distributed SVD.
- 
- 4.
- 5. **Randomized PCA / Approximate PCA:**
  - **Concept:** Instead of computing the exact PCA, use a randomized algorithm to find a very good approximation much faster.
  - **Method:** Randomized SVD, as discussed previously, can be highly parallelized. The main step involves computing  $Y = XR$ , where R is a random matrix. This matrix multiplication is easily parallelizable.
  - **Benefit:** This is often the most practical and efficient approach for very large matrices, as implemented in frameworks like Spark.
- 6.

**Conclusion:** For massive datasets, parallel and approximate methods are essential. Apache Spark's MLlib provides robust, off-the-shelf implementations of PCA that use these strategies to scale to terabytes of data.

---

## Question 8

**Compare the use of PCA to select features with other feature selection methods.**

### Theory

Using PCA for feature engineering is fundamentally different from traditional feature selection methods. PCA **transforms** features, while feature selection **selects** a subset of the original features.

### PCA for Feature Engineering

- **Method:** PCA creates a new set of features (the principal components) that are linear combinations of *all* the original features. It does not select from the original features; it creates new ones.
- **Pros:**
  - It can capture information from all original variables, even if in a compressed form.
  - It produces uncorrelated features, which can be beneficial for some models.
- 
- **Cons:**

- **Loss of Interpretability:** The new principal components are abstract mathematical constructs. It is often difficult to relate them back to the original, meaningful business or scientific features. A component might be  $0.6*\text{age} - 0.3*\text{income} + 0.7*$ ..., which is not directly interpretable.
- 

## Traditional Feature Selection Methods

These methods select a subset of the original features and discard the rest.

1. **Filter Methods:**
  - **Method:** Rank features based on a statistical metric (e.g., correlation with the target, mutual information, variance) and select the top k features.
  - **Pros:** Fast and model-agnostic.
  - **Cons:** Ignores feature interactions.
- 2.
3. **Wrapper Methods:**
  - **Method:** Use a machine learning model to evaluate different subsets of features and select the subset that gives the best model performance (e.g., Recursive Feature Elimination).
  - **Pros:** Can find the optimal feature set for a specific model.
  - **Cons:** Computationally very expensive.
- 4.
5. **Embedded Methods:**
  - **Method:** Perform feature selection as part of the model training process (e.g., L1 regularization in Lasso, which drives the coefficients of unimportant features to zero).
  - **Pros:** Efficient and considers feature interactions in the context of the model.
- 6.

## Comparison and When to Use Which

Feature	PCA (Feature Transformation)	Feature Selection
Output	A new, smaller set of synthetic features.	A subset of the original features.
Interpretability	<b>Low.</b> The new features are abstract.	<b>High.</b> The selected features are the original, interpretable ones.
Goal	Maximize variance preservation.	Varies (e.g., maximize model performance, reduce redundancy).
Information Loss	Can retain information from all original features.	Information from discarded features is completely lost.

- **Use Feature Selection when:**
    - **Interpretability is a priority.** You need to be able to explain the model's predictions in terms of the original variables.
    - You want to reduce the cost of data collection by identifying a smaller set of necessary features.
  - **Use PCA when:**
    - **Interpretability is less important** than predictive performance.
    - You have a large number of highly correlated features and want to handle multicollinearity.
    - Your primary goal is dimensionality reduction to speed up a model, and you are willing to trade interpretability for performance.
  -
- 

## Question 9

**What recent advancements have been made concerning PCA for big data?**

### Theory

The primary advancements concerning PCA for big data have focused on overcoming the computational and memory bottlenecks of standard PCA. These advancements fall into two main categories: **randomized/approximate algorithms** and **distributed/parallel computing frameworks**.

### Key Advancements

1. **Randomized PCA and Randomized SVD:**
  - **Advancement:** This is perhaps the most significant practical advancement. Instead of computing the full, exact SVD, randomized algorithms use probabilistic methods to find a low-rank approximation of the data matrix very quickly.
  - **Benefit:** They are much faster and more memory-efficient than traditional methods, especially when you only need the top few principal components. The approximation is typically very close to the exact result.
  - **Impact:** This has made PCA feasible for very large and high-dimensional matrices on a single machine and is the core of modern, efficient implementations (e.g., in scikit-learn).
- 2.
3. **Distributed and Parallel PCA on Frameworks like Apache Spark:**
  - **Advancement:** Algorithms have been developed to perform PCA on data that is distributed across a cluster of machines.

- **Method:** Spark's MLlib, for instance, implements several strategies. It can compute the covariance matrix in a distributed manner or use distributed SVD algorithms (like block-based methods or TSQR) to decompose the data matrix directly.
  - **Impact:** This allows PCA to scale to massive, terabyte-scale datasets that would be impossible to process on a single computer.
- 4.
5. **Streaming and Online PCA:**
- **Advancement:** Development of algorithms, like **Incremental PCA**, that can update a PCA model as new data arrives in a stream, without needing to reprocess all historical data.
  - **Impact:** This enables the use of PCA in real-time applications for monitoring, anomaly detection, and adaptive modeling.
- 6.
7. **GPU Acceleration:**
- **Advancement:** Leveraging GPUs to accelerate the linear algebra operations (matrix multiplications) at the heart of PCA and SVD.
  - **Impact:** For moderately sized datasets that fit on a GPU's memory, this can provide a massive speedup over CPU-based computations. Libraries like cuML from RAPIDS provide GPU-accelerated versions of PCA.
- 8.

These advancements have transformed PCA from a technique limited to moderate-sized datasets into a scalable, cornerstone algorithm for modern big data analytics.

---

## Pca Interview Questions - Coding Questions

### Question 1

**Write a Python function to perform PCA from scratch using NumPy.**

#### Theory

PCA from scratch involves a sequence of clear mathematical steps: standardizing the data, calculating the covariance matrix, performing an eigen-decomposition on this matrix to find the principal components (eigenvectors) and their corresponding variances (eigenvalues), and finally projecting the data onto the selected components.

#### Code Example

```

code Python
downloadcontent_copyexpand_less
import numpy as np

def pca_from_scratch(X, n_components):
    """
    Performs PCA from scratch using NumPy.

    Args:
        X (np.array): The input data, shape (n_samples, n_features).
        n_components (int): The number of principal components to keep.

    Returns:
        np.array: The transformed data, shape (n_samples, n_components).
    """

    # 1. Standardize the data (center and scale)
    # Center the data by subtracting the mean
    X_centered = X - np.mean(X, axis=0)

    # 2. Calculate the covariance matrix
    # (n_samples-1) is used for the sample covariance
    cov_matrix = np.cov(X_centered, rowvar=False)

    # 3. Perform Eigen-decomposition on the covariance matrix
    eigenvalues, eigenvectors = np.linalg.eig(cov_matrix)

    # 4. Sort eigenvectors by descending eigenvalues
    # Create pairs of (eigenvalue, eigenvector) and sort
    eigen_pairs = [(np.abs(eigenvalues[i]), eigenvectors[:, i]) for i in range(len(eigenvalues))]
    eigen_pairs.sort(key=lambda k: k[0], reverse=True)

    # 5. Select the top k eigenvectors (principal components)
    top_eigenvectors = np.array([pair[1] for pair in eigen_pairs[:n_components]]).T

    # 6. Project the original data onto the principal components
    X_projected = X_centered.dot(top_eigenvectors)

    return X_projected

# --- Example Usage ---
if __name__ == '__main__':
    from sklearn.datasets import load_iris
    import matplotlib.pyplot as plt

```

```

# Load Iris dataset
iris = load_iris()
X = iris.data
y = iris.target

# Apply PCA to reduce to 2 components
X_pca = pca_from_scratch(X, 2)

# Plot the results
plt.figure(figsize=(8, 6))
plt.scatter(X_pca[:, 0], X_pca[:, 1], c=y, cmap='viridis')
plt.xlabel('First Principal Component')
plt.ylabel('Second Principal Component')
plt.title('PCA from Scratch on Iris Dataset')
plt.grid(True)
plt.show()

```

## Explanation

- Standardize Data:** We first center the data by subtracting the mean of each feature (column). `np.mean(X, axis=0)` calculates the mean for each column.
  - Covariance Matrix:** `np.cov(X_centered, rowvar=False)` computes the covariance matrix. `rowvar=False` is important; it tells NumPy that columns represent variables and rows represent observations.
  - Eigen-decomposition:** `np.linalg.eig()` performs the eigen-decomposition, returning the eigenvalues and their corresponding eigenvectors.
  - Sort Components:** We sort the eigenvectors based on the magnitude of their corresponding eigenvalues in descending order. This ranks the principal components by the amount of variance they explain.
  - Select Components:** We construct a projection matrix from the top `n_components` eigenvectors.
  - Project Data:** Finally, we project the centered data onto these top eigenvectors using a dot product to get the new, lower-dimensional representation.
- 

## Question 2

**Use scikit-learn to apply PCA on a high-dimensional dataset and interpret the results.**

### Theory

scikit-learn provides a clean, easy-to-use implementation of PCA. The process involves standardizing the data, fitting the PCA object, and then transforming the data. The results can be interpreted by examining the `explained_variance_ratio_` and the `components_` (loadings) attributes of the fitted object.

### Code Example

```
code Python
downloadcontent_copyexpand_less
IGNORE_WHEN COPYING_START
IGNORE_WHEN COPYING_END
    import numpy as np
import matplotlib.pyplot as plt
from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler
from sklearn.datasets import fetch_olivetti_faces

# 1. Load a high-dimensional dataset
# The Olivetti faces dataset has 400 images, each 64x64 pixels (4096 features).
faces = fetch_olivetti_faces()
X = faces.data
n_samples, n_features = X.shape
print(f"Dataset shape: {n_samples} samples, {n_features} features")

# 2. Standardize the data
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# 3. Apply PCA
# Let's find the number of components to explain 95% of the variance
pca = PCA(n_components=0.95, svd_solver='full')
pca.fit(X_scaled)

# 4. Interpret the results
# How many components were chosen?
n_components_chosen = pca.n_components_
print(f"\nNumber of components to explain 95% variance: {n_components_chosen}")

# Plot the cumulative explained variance
plt.figure(figsize=(10, 6))
plt.plot(np.cumsum(pca.explained_variance_ratio_), marker='o', linestyle='--')
plt.xlabel('Number of Components')
plt.ylabel('Cumulative Explained Variance')
plt.title('Explained Variance by Number of Components')
plt.axhline(y=0.95, color='r', linestyle='-', label='95% Threshold')
```

```

plt.grid(True)
plt.legend()
plt.show()

# Visualize the first few principal components (eigenfaces)
fig, axes = plt.subplots(2, 5, figsize=(12, 5),
                       subplot_kw={'xticks':[], 'yticks':[]},
                       gridspec_kw=dict(hspace=0.1, wspace=0.1))

for i, ax in enumerate(axes.flat):
    ax.imshow(pca.components_[i].reshape(64, 64), cmap='bone')
    ax.set_title(f"PC {i+1}")

plt.suptitle("Top 10 Principal Components (Eigenfaces)")
plt.show()

```

## Explanation and Interpretation

1. **Load and Scale Data:** We load the Olivetti faces dataset, which is high-dimensional (4096 features). We standardize it as a best practice.
  2. **Apply PCA:** We instantiate PCA with n\_components=0.95. This tells scikit-learn to automatically select the minimum number of components required to preserve 95% of the total variance in the data.
  3. **Interpret Explained Variance:**
    - The output shows that we reduced the dimensionality from 4096 features down to a much smaller number (e.g., around 123) while still keeping 95% of the information.
    - The cumulative explained variance plot confirms this, showing a steep initial rise (the first few components capture a lot of variance) and then leveling off. The red line indicates our 95% threshold.
  - 4.
  5. **Interpret Components (Eigenfaces):**
    - We access the principal components via pca.components\_.
    - By reshaping each component back into a 64x64 image and plotting it, we can visualize the "eigenfaces." These are the fundamental patterns of variation in the dataset. PC1 might represent the average face, while subsequent components might represent variations like lighting from the side, the presence of glasses, or different facial expressions.
  - 6.
- 

## Question 3

**Code a Python script to visualize the eigenfaces from a given set of facial images dataset using PCA.**

## Theory

The "Eigenfaces" method is a classic application of PCA for facial recognition and analysis. It treats each face image as a high-dimensional vector and applies PCA to find the principal components of the face dataset. These components, when visualized as images, are called eigenfaces because they are the eigenvectors of the covariance matrix of the face images. They represent the fundamental "features" or patterns that constitute a face.

## Code Example

```
code Python
downloadcontent_copyexpand_less
IGNORE_WHEN COPYING_START
IGNORE_WHEN COPYING_END
    import matplotlib.pyplot as plt
from sklearn.datasets import fetch_lfw_people
from sklearn.decomposition import PCA

# 1. Load the Labeled Faces in the Wild (LFW) dataset
# min_faces_per_person filters for people with at least 70 images
# This provides a sizable dataset for PCA.
lfw_people = fetch_lfw_people(min_faces_per_person=70, resize=0.4)
X = lfw_people.data
n_samples, h, w = lfw_people.images.shape
print(f"Dataset loaded: {n_samples} faces, each {h}x{w} pixels")

# 2. Apply PCA to extract the eigenfaces
# We'll extract the top 150 components
n_components = 150
pca = PCA(n_components=n_components, svd_solver='randomized', whiten=True)
pca.fit(X)

# The principal components are the eigenfaces
eigenfaces = pca.components_.reshape((n_components, h, w))

# 3. Visualize the top eigenfaces
def plot_gallery(images, titles, h, w, n_row=3, n_col=4):
    """Helper function to plot a gallery of portraits"""
    plt.figure(figsize=(1.8 * n_col, 2.4 * n_row))
    plt.subplots_adjust(bottom=0, left=.01, right=.99, top=.90, hspace=.35)
    for i in range(n_row * n_col):
        plt.subplot(n_row, n_col, i + 1)
        plt.imshow(images[i], cmap=plt.cm.gray)
        if titles[i] != '':
            plt.title(titles[i])
```

```

plt.title(titles[i], size=12)
plt.xticks(())
plt.yticks(())

eigenface_titles = [f"Eigenface {i+1}" for i in range(eigenfaces.shape[0])]
plot_gallery(eigenfaces, eigenface_titles, h, w)
plt.suptitle("Top 12 Eigenfaces", size=16)
plt.show()

```

## Explanation

1. **Load Data:** We use the fetch\_lfw\_people dataset from scikit-learn, which contains a large number of facial images. Each image is flattened into a vector of pixel intensities.
  2. **Apply PCA:** We instantiate PCA to extract n\_components=150.
    - svd\_solver='randomized': This is used for efficiency on high-dimensional data where we only want the top components.
    - whiten=True: This step scales the principal components to have unit variance. While not always necessary, it often improves the visual contrast of the eigenfaces, making the patterns more distinct.
  - 3.
  4. **Extract Eigenfaces:** The principal components are stored in the pca.components\_ attribute. Each row of this matrix is an eigenvector. We reshape each row back into the original image dimensions (h x w) to get the visual representation of the eigenface.
  5. **Visualize:** We use a helper function plot\_gallery to display the first 12 eigenfaces in a grid. The resulting images are often ghost-like and highlight the most significant variations in the face dataset, such as lighting, hairstyle, and general facial structure.
- 

## Question 4

**Implement PCA for feature extraction before applying a machine learning model.**

### Theory

PCA can be used as a powerful feature extraction technique. Instead of training a machine learning model on the original high-dimensional and potentially correlated features, we can train it on the lower-dimensional, uncorrelated principal components. This can lead to faster training times, reduced overfitting, and sometimes even better performance by filtering out noise.

### Code Example

code Python  
[downloadcontent\\_copyexpand\\_less](#)

```

IGNORE_WHEN COPYING_START
IGNORE_WHEN COPYING_END
    import numpy as np
from sklearn.datasets import load_digits
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
import time

# 1. Load the dataset
digits = load_digits()
X, y = digits.data, digits.target
print(f"Original data shape: {X.shape}") # (1797 samples, 64 features)

# 2. Split data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# --- Scenario 1: Model without PCA ---
print("\n--- Training Model on Original Data ---")
scaler_orig = StandardScaler()
X_train_orig_scaled = scaler_orig.fit_transform(X_train)
X_test_orig_scaled = scaler_orig.transform(X_test)

logreg_orig = LogisticRegression(max_iter=1000)
start_time = time.time()
logreg_orig.fit(X_train_orig_scaled, y_train)
end_time = time.time()

preds_orig = logreg_orig.predict(X_test_orig_scaled)
acc_orig = accuracy_score(y_test, preds_orig)
print(f"Accuracy on original 64 features: {acc_orig:.4f}")
print(f"Training time: {end_time - start_time:.4f} seconds")

# --- Scenario 2: Model with PCA ---
print("\n--- Training Model on PCA-transformed Data ---")
scaler_pca = StandardScaler()
X_train_scaled = scaler_pca.fit_transform(X_train)
X_test_scaled = scaler_pca.transform(X_test)

# Apply PCA to retain 95% of variance
pca = PCA(n_components=0.95)
X_train_pca = pca.fit_transform(X_train_scaled)

```

```

X_test_pca = pca.transform(X_test_scaled)
print(f"PCA-transformed data shape: {X_train_pca.shape}")

logreg_pca = LogisticRegression(max_iter=1000)
start_time = time.time()
logreg_pca.fit(X_train_pca, y_train)
end_time = time.time()

preds_pca = logreg_pca.predict(X_test_pca)
acc_pca = accuracy_score(y_test, preds_pca)
print(f"Accuracy on {pca.n_components_} principal components: {acc_pca:.4f}")
print(f"Training time: {end_time - start_time:.4f} seconds")

```

## Explanation

1. **Setup:** We load the digits dataset, which has 64 features (8x8 pixel images), and split it into training and testing sets.
  2. **Model without PCA:** As a baseline, we train a LogisticRegression model on the original (but scaled) 64 features. We measure its accuracy and training time.
  3. **Model with PCA:**
    - We apply StandardScaler and then PCA to the training data. We set n\_components=0.95 to automatically select enough components to keep 95% of the variance.
    - We use the **same fitted pca object** to transform the test data. This is a critical step to ensure consistency.
    - The dimensionality is significantly reduced (e.g., from 64 to around 29 features).
    - We then train a new LogisticRegression model on this lower-dimensional X\_train\_pca.
  - 4.
  5. **Comparison:** We compare the two scenarios. We typically observe:
    - **Similar Accuracy:** The accuracy of the PCA-based model is often very close to the original, and sometimes slightly better due to noise reduction.
    - **Faster Training:** The training time for the PCA-based model is significantly lower because it is working with fewer features. This benefit becomes much more pronounced with more complex models or larger datasets.
  - 6.
- 

## Question 5

**Implement PCA in TensorFlow or PyTorch and compare the results with scikit-learn's implementation.**

## Theory

PCA can be implemented in deep learning frameworks like TensorFlow or PyTorch by leveraging their efficient tensor operations to perform the required linear algebra: centering the data, computing the covariance matrix, and finding the eigenvalues and eigenvectors. The results should be mathematically equivalent to scikit-learn's implementation, though minor differences can occur due to floating-point precision.

### Code Example (Using TensorFlow)

```
code Python
downloadcontent_copyexpand_less
IGNORE_WHEN COPYING_START
IGNORE_WHEN COPYING_END
    import numpy as np
    import tensorflow as tf
    from sklearn.decomposition import PCA as SklearnPCA
    from sklearn.preprocessing import StandardScaler
    from sklearn.datasets import make_blobs

def pca_tensorflow(X_tensor, n_components):
    """Performs PCA using TensorFlow."""
    # 1. Center the data
    mean = tf.reduce_mean(X_tensor, axis=0)
    X_centered = X_tensor - mean

    # 2. Calculate the covariance matrix
    n_samples = tf.cast(tf.shape(X_tensor)[0], tf.float32)
    # Note: tf.linalg.eigh requires a symmetric matrix
    cov_matrix = tf.matmul(tf.transpose(X_centered), X_centered) / (n_samples - 1)

    # 3. Perform Eigen-decomposition
    # tf.linalg.eigh is used for symmetric/Hermitian matrices
    eigenvalues, eigenvectors = tf.linalg.eigh(cov_matrix)

    # 4. Sort eigenvectors by descending eigenvalues
    # TensorFlow returns eigenvalues in ascending order, so we reverse them
    sorted_indices = tf.argsort(eigenvalues, direction='DESCENDING')
    sorted_eigenvectors = tf.gather(eigenvectors, sorted_indices, axis=1)

    # 5. Select top k eigenvectors
    top_eigenvectors = sorted_eigenvectors[:, :n_components]

    # 6. Project the data
    X_projected = tf.matmul(X_centered, top_eigenvectors)
```

```

return X_projected

# --- Example Usage ---
if __name__ == '__main__':
    # Generate sample data
    X, y = make_blobs(n_samples=1000, n_features=10, centers=5, random_state=42)
    X_scaled = StandardScaler().fit_transform(X)
    n_components = 2

    # --- TensorFlow PCA ---
    print("--- TensorFlow PCA ---")
    X_tf_tensor = tf.convert_to_tensor(X_scaled, dtype=tf.float32)
    X_pca_tf = pca_tensorflow(X_tf_tensor, n_components)
    print("First 5 projected samples (TF):\n", X_pca_tf.numpy()[:5])

    # --- Scikit-learn PCA ---
    print("\n--- Scikit-learn PCA ---")
    sklearn_pca = SklearnPCA(n_components=n_components)
    X_pca_sklearn = sklearn_pca.fit_transform(X_scaled)
    print("First 5 projected samples (sklearn):\n", X_pca_sklearn[:5])

# Note: The signs of the eigenvectors can be flipped (-v is also a valid eigenvector).
# The resulting projections might be reflections of each other, but the structure is identical.
# We can check this by comparing the absolute values.
abs_diff = np.mean(np.abs(np.abs(X_pca_tf.numpy()) - np.abs(X_pca_sklearn)))
print(f"\nMean absolute difference between projections: {abs_diff:.6f}")

```

## Explanation

- TensorFlow Implementation:** The `pca_tensorflow` function mirrors the from-scratch NumPy implementation but uses TensorFlow's tensor operations (`tf.reduce_mean`, `tf.matmul`, `tf.linalg.eigh`). This allows the computation to potentially run on a GPU.
- Scikit-learn Comparison:** We run scikit-learn's PCA on the same scaled data for a direct comparison.
- Result Analysis:**
  - The numerical results from both implementations are virtually identical.
  - An important point to note is that the **sign of eigenvectors is arbitrary**. An eigenvector  $v$  and its negative  $-v$  are both valid eigenvectors. Because of this, the output of the TensorFlow implementation might be a mirror image of the scikit-learn output (e.g., all signs flipped in one component).
  - This sign flip is not an error. The geometric structure, distances, and variance explained are identical. We verify this by comparing the absolute values of the projected data, which shows a near-zero difference.

---

## Question 6

**Demonstrate how to choose an optimal number of dimensions with PCA in Python using the “elbow method”.**

### Theory

The "elbow method" in the context of PCA refers to plotting the explained variance of each principal component to find a point where the variance explained starts to level off. This "elbow" suggests a natural cut-off point for the number of components to retain, as subsequent components contribute diminishing amounts of information.

### Code Example

```
code Python
downloadcontent_copyexpand_less
IGNORE_WHEN COPYING_START
IGNORE_WHEN COPYING_END
    import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import make_classification
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA

# 1. Generate a synthetic high-dimensional dataset
X, y = make_classification(
    n_samples=1000,
    n_features=50,
    n_informative=15,
    n_redundant=10,
    random_state=42
)

# 2. Standardize the data
X_scaled = StandardScaler().fit_transform(X)

# 3. Apply PCA to get all possible components
# By not setting n_components, we get all 50 components
pca = PCA()
pca.fit(X_scaled)
```

```

# 4. Create the scree plot (the "elbow method" plot)
explained_variance = pca.explained_variance_
n_components = np.arange(1, len(explained_variance) + 1)

plt.figure(figsize=(12, 7))
plt.plot(n_components, explained_variance, 'o-', linewidth=2, color='blue')
plt.title('Scree Plot (Elbow Method)')
plt.xlabel('Principal Component')
plt.ylabel('Eigenvalue (Explained Variance)')
plt.grid(True)

# Annotate a potential "elbow"
# This part is subjective, based on visual inspection
elbow_point = 15
plt.axvline(x=elbow_point, color='r', linestyle='--', label=f'Potential Elbow at k={elbow_point}')
plt.legend()

plt.show()

# --- Interpretation ---
# By examining the plot, we look for the "knee" or "elbow" of the curve.
# In this example, the curve starts to flatten out significantly after about the 15th component.
# This suggests that the first ~15 components capture the most significant variance (the
# "signal"),
# while the subsequent components capture less information and may represent noise.
# Therefore, a good choice for the number of dimensions to keep would be around 15.

```

## Explanation

- Generate Data:** We create a synthetic dataset with 50 features, but only 15 of them are truly informative, giving us a known underlying dimensionality.
  - Fit PCA:** We fit PCA without specifying n\_components. This computes all 50 principal components and their corresponding explained variances (eigenvalues).
  - Plot the Scree Plot:** We plot the component number against its explained variance. This plot is often called a **scree plot**.
  - Identify the Elbow:** We visually inspect the plot to find the "elbow." In the generated plot, there is a clear drop-off that begins to level out around k=15. This aligns with the n\_informative=15 parameter we used to generate the data. The red dashed line highlights this subjective choice.
  - Conclusion:** Based on this visual heuristic, we would conclude that keeping approximately 15 components is a good choice, as it balances dimensionality reduction with information preservation.
-

## Question 7

**Write a Python script to automatically remove outliers before performing PCA.**

### Theory

Standard PCA is sensitive to outliers because it is based on variance, which is itself sensitive to extreme values. Outliers can heavily influence the calculation of the mean and covariance matrix, thus skewing the principal components. A common approach to make PCA more robust is to first identify and remove outliers from the dataset. One simple yet effective method for outlier detection is using the Mahalanobis distance or an isolation forest.

Here, we'll use the **Isolation Forest** algorithm, which is efficient and works well in high-dimensional spaces.

### Code Example

```
code Python
downloadcontent_copyexpand_less
IGNORE_WHEN COPYING_START
IGNORE_WHEN COPYING_END

import numpy as np
from sklearn.ensemble import IsolationForest
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
import matplotlib.pyplot as plt

# 1. Create a synthetic dataset with some outliers
np.random.seed(42)
X_inliers = 0.5 * np.random.randn(500, 2)
X_outliers = np.random.uniform(low=-4, high=4, size=(20, 2))
X = np.vstack([X_inliers, X_outliers])

# 2. Identify outliers using Isolation Forest
# contamination='auto' is a good default, or you can specify a fraction
iso_forest = IsolationForest(contamination='auto', random_state=42)
outlier_preds = iso_forest.fit_predict(X)

# The algorithm returns -1 for outliers and 1 for inliers
outliers_mask = outlier_preds == -1
inliers_mask = outlier_preds == 1

X_inliers_data = X[inliers_mask, :]
X_outliers_data = X[outliers_mask, :]
print(f"Original dataset size: {len(X)}")
```

```

print(f"Number of outliers detected: {len(X_outliers_data)}")

# 3. Perform PCA on both the original and cleaned data for comparison
# Standardize data first
scaler = StandardScaler()

# PCA on original data
X_scaled_orig = scaler.fit_transform(X)
pca_orig = PCA(n_components=2)
pca_orig.fit(X_scaled_orig)
pc1_orig = pca_orig.components_[0]

# PCA on cleaned data (inliers only)
X_scaled_clean = scaler.fit_transform(X_inliers_data)
pca_clean = PCA(n_components=2)
pca_clean.fit(X_scaled_clean)
pc1_clean = pca_clean.components_[0]

# 4. Visualize the results
plt.figure(figsize=(14, 6))

# Plot 1: Original data with outliers and original PC1
plt.subplot(1, 2, 1)
plt.scatter(X_inliers_data[:, 0], X_inliers_data[:, 1], label='Inliers')
plt.scatter(X_outliers_data[:, 0], X_outliers_data[:, 1], color='r', label='Outliers')
# Plot the PC1 vector (scaled for visibility)
plt.quiver(0, 0, pc1_orig[0]*3, pc1_orig[1]*3, angles='xy', scale_units='xy', scale=1, color='black',
label='Original PC1')
plt.title('PCA on Original Data (with Outliers)')
plt.legend()
plt.axis('equal')

# Plot 2: Cleaned data and new PC1
plt.subplot(1, 2, 2)
plt.scatter(X_inliers_data[:, 0], X_inliers_data[:, 1], label='Inliers')
# Plot the PC1 vector from the cleaned data
plt.quiver(0, 0, pc1_clean[0]*3, pc1_clean[1]*3, angles='xy', scale_units='xy', scale=1,
color='green', label='Cleaned PC1')
plt.title('PCA on Cleaned Data (Outliers Removed)')
plt.legend()
plt.axis('equal')

plt.show()

```

## Explanation

1. **Generate Data with Outliers:** We create a dataset with a clear cluster of "inlier" points and some scattered "outlier" points.
2. **Detect Outliers:** We use IsolationForest, an unsupervised anomaly detection algorithm. It works by isolating observations and returns -1 for anomalies (outliers) and 1 for normal points (inliers). We use this to create a mask to separate the two groups.
3. **Perform PCA:** We run PCA twice: once on the full, original dataset and once on the "cleaned" dataset containing only the inliers.
4. **Visualize and Compare:** The visualization clearly shows the effect.
  - o In the left plot, the **original PC1** (black arrow) is skewed by the outliers. It is pulled away from the main direction of variance of the inlier data cloud.
  - o In the right plot, after removing the outliers, the **cleaned PC1** (green arrow) correctly aligns with the true direction of maximum variance of the main data cluster.
- 5.

This demonstrates that removing outliers first leads to a more accurate and robust estimation of the principal components.

---

## Question 8

**Create a synthetic dataset and show the effect of PCA on classification accuracy using a machine learning algorithm before and after PCA.**

### Theory

Applying PCA before a classification algorithm can have several effects. It can:

1. **Reduce Overfitting:** By reducing the number of features, especially noisy or redundant ones.
2. **Speed up Training:** Fewer features lead to faster model training.
3. **Maintain or Improve Accuracy:** If the principal components capture the signal well, accuracy can remain high or even improve due to noise reduction.
4. **Decrease Accuracy:** If the components with low variance are important for classification, PCA might discard them and hurt performance.

This script demonstrates a common scenario where PCA helps by removing redundant features.

### Code Example

```
code Python  
downloadcontent_copyexpand_less  
IGNORE_WHEN COPYING_START
```

```

IGNORE_WHEN COPYING-END
import numpy as np
from sklearn.datasets import make_classification
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score

# 1. Create a synthetic dataset
# High number of features, with many being redundant.
X, y = make_classification(
    n_samples=1000,
    n_features=100,
    n_informative=10,
    n_redundant=80, # 80 redundant features
    n_classes=2,
    random_state=42
)
print(f"Original feature count: {X.shape[1]}")

# Split into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# --- Scenario 1: SVM on Original Features ---
scaler_orig = StandardScaler()
X_train_orig = scaler_orig.fit_transform(X_train)
X_test_orig = scaler_orig.transform(X_test)

svm_orig = SVC(kernel='linear', random_state=42)
svm_orig.fit(X_train_orig, y_train)
y_pred_orig = svm_orig.predict(X_test_orig)
acc_orig = accuracy_score(y_test, y_pred_orig)

print(f"\nAccuracy with original 100 features: {acc_orig:.4f}")

# --- Scenario 2: SVM on PCA-transformed Features ---
scaler_pca = StandardScaler()
X_train_pca_prep = scaler_pca.fit_transform(X_train)
X_test_pca_prep = scaler_pca.transform(X_test)

# Let's keep enough components to explain 95% of the variance
pca = PCA(n_components=0.95)
X_train_pca = pca.fit_transform(X_train_pca_prep)

```

```

X_test_pca = pca.transform(X_test_pca_prep)
print(f"Number of features after PCA: {X_train_pca.shape[1]}")

svm_pca = SVC(kernel='linear', random_state=42)
svm_pca.fit(X_train_pca, y_train)
y_pred_pca = svm_pca.predict(X_test_pca)
acc_pca = accuracy_score(y_test, y_pred_pca)

print(f"Accuracy with {pca.n_components_} principal components: {acc_pca:.4f}")

# --- Conclusion ---
# In this case, the accuracy is often very similar or even slightly better after PCA.
# This is because PCA effectively removed the redundant information, simplifying
# the problem for the SVM classifier without losing the important signal.
# A significant benefit would also be a reduction in training time.

```

## Explanation

- Generate Data:** We use `make_classification` to create a dataset with 100 features, but deliberately make 80 of them redundant (linear combinations of the 10 informative features). This simulates a common real-world scenario.
  - Baseline Model:** We train a Support Vector Machine (SVM) classifier on all 100 original (scaled) features and record its accuracy on the test set.
  - PCA Model:** We create a PCA pipeline. We scale the data and then apply `PCA(n_components=0.95)`. PCA will discover that most of the variance is contained in a much smaller number of components because of the high redundancy, and it will reduce the dimensionality accordingly.
  - Train on Components:** We train a new SVM on the resulting low-dimensional principal components.
  - Compare Results:** We compare the accuracy scores. In this typical example, the accuracy after PCA is very close to the original, demonstrating that PCA successfully reduced the feature space from 100 to a much smaller number (e.g., ~15-20) without losing predictive power. This makes the model more efficient and potentially more robust against overfitting.
- 

# Pca Interview Questions - Scenario\_Based Questions

# Question 1

Discuss the importance of the trace of a matrix in the context of PCA.

## Theory

In the context of PCA, the **trace** of a matrix (the sum of the elements on its main diagonal) is important because the **trace of the covariance matrix is equal to the total variance of the dataset**. This provides a fundamental link between the original data's variance and the eigenvalues derived from PCA.

## Explanation

### 1. Trace of the Covariance Matrix:

- Let  $C$  be the  $p \times p$  covariance matrix of a dataset with  $p$  features.
- The diagonal elements of  $C$  ( $C_{11}, C_{22}, \dots, C_{pp}$ ) are the variances of each individual feature ( $\text{Var}(X_1), \text{Var}(X_2), \dots, \text{Var}(X_p)$ ).
- The trace of  $C$  is  $\text{Tr}(C) = \sum_i C_{ii} = \sum_i \text{Var}(X_i)$ .
- Therefore,  $\text{Tr}(C)$  represents the **total variance** of the dataset.

2.

### 3. Trace and Eigenvalues:

- A fundamental property of matrices is that the trace of a matrix is equal to the sum of its eigenvalues.
- Let  $\lambda_1, \lambda_2, \dots, \lambda_p$  be the eigenvalues of the covariance matrix  $C$ .
- Then,  $\text{Tr}(C) = \sum_i \lambda_i$ .

4.

### 5. The Connection to PCA:

- From point 1, we know  $\text{Tr}(C) = \text{Total Variance}$ .
- From point 2, we know  $\text{Tr}(C) = \text{Sum of Eigenvalues}$ .
- Therefore, **Total Variance = Sum of Eigenvalues**.
- In PCA, each eigenvalue  $\lambda_i$  represents the variance explained by its corresponding principal component. This property confirms that the sum of the variances of all principal components is equal to the total variance of the original dataset.

6.

## Why This is Important

This property provides the mathematical foundation for calculating the **explained variance ratio**.

The proportion of total variance explained by a single principal component  $PC_i$  is:

$$\text{Explained Variance Ratio} = \lambda_i / \text{Total Variance} = \lambda_i / \text{Tr}(C) = \lambda_i / \sum \lambda_j$$

This allows us to quantitatively assess how much information is preserved when we reduce the dimensionality of the data, making the entire process interpretable and controllable.

---

## Question 2

**Discuss the application of PCA in feature engineering.**

### Theory

PCA is a powerful technique in feature engineering, where the goal is to create new features from existing ones to improve machine learning model performance. PCA contributes by transforming a set of potentially correlated, high-dimensional features into a new set of uncorrelated, lower-dimensional features called principal components. These new features can then be used to train a model.

### Applications in Feature Engineering

#### 1. Dimensionality Reduction:

- **Application:** When a dataset has a very large number of features (high dimensionality), models can become slow to train and are prone to overfitting.
- **PCA's Role:** PCA can reduce the number of features to a more manageable size while retaining most of the original variance. This creates a more compact and efficient feature set.

#### 2.

#### 3. Handling Multicollinearity:

- **Application:** Multicollinearity occurs when features are highly correlated with each other. This can cause problems for certain models, particularly linear models like Linear Regression, by making their coefficient estimates unstable and difficult to interpret.
- **PCA's Role:** By design, the principal components produced by PCA are **orthogonal** (perpendicular) and therefore **uncorrelated**. By using the principal components as features instead of the original variables, the problem of multicollinearity is completely eliminated.

#### 4.

#### 5. Noise Reduction:

- **Application:** Real-world data is often noisy. This noise can degrade model performance.
- **PCA's Role:** PCA assumes that the signal is concentrated in the components with high variance, while noise is in the components with low variance. By selecting only the top k components, PCA acts as a denoising filter, creating a new feature set that is potentially cleaner and leads to a more robust model.

#### 6.

### Trade-offs

The main trade-off when using PCA for feature engineering is the **loss of interpretability**.

- The original features (e.g., age, income) have a clear meaning.
- The new features created by PCA are abstract linear combinations of all original features (e.g.,  $PC1 = 0.7*age - 0.4*income + \dots$ ). While they are mathematically optimal for capturing variance, they are difficult to explain in a business context.

**Conclusion:** PCA is an excellent feature engineering tool when the primary goal is to improve the predictive performance, speed, or stability of a model, and when direct interpretability of the features is a lower priority.

---

## Question 3

**Discuss how PCA can suffer from outlier sensitivity and ways to address it.**

### Theory

Standard PCA is **highly sensitive to outliers**. This is a direct consequence of its objective function, which is to maximize variance (or minimize the sum of squared reconstruction errors). Outliers, being extreme values, can disproportionately inflate the variance in their direction, thus heavily influencing the calculation of the principal components.

### How Outliers Affect PCA

1. **Skewing the Mean:** PCA begins by centering the data. Outliers will pull the mean of a feature towards them, causing the centered data to be misrepresented.
2. **Distorting the Covariance Matrix:** The variance and covariance calculations are based on squared distances. An outlier far from the mean will have a very large squared distance, giving it a huge influence on the variance calculation for its feature and its covariance with other features.
3. **Hijacking Principal Components:** Because PCA seeks the direction of maximum variance, a single or a few outliers can create a spurious direction of high variance. The first principal component, instead of capturing the main structure of the bulk of the data, may simply point from the data's center towards the outlier(s). This renders the component useless for describing the true data structure.

### Ways to Address Outlier Sensitivity

1. **Outlier Detection and Removal (Preprocessing):**
  - **Method:** Before applying PCA, use a robust outlier detection algorithm to identify and remove anomalous data points.
  - **Techniques:**
    - **Isolation Forest:** An efficient algorithm that works well in high dimensions.

- **Local Outlier Factor (LOF)**: A density-based method that identifies outliers based on their local neighborhood.
  - **Elliptic Envelope**: Assumes the inlier data follows a Gaussian distribution and identifies points outside this distribution.
  - 
  - **Pros/Cons**: This is a direct and effective approach, but it involves removing data, which may not always be desirable.
- 2.
3. **Use of Robust Scaling:**
- **Method**: Instead of using StandardScaler (which is based on mean and standard deviation and is sensitive to outliers), use a RobustScaler.
  - **Technique**: RobustScaler centers the data using the **median** and scales it using the **interquartile range (IQR)**. Both median and IQR are much less sensitive to outliers than the mean and standard deviation.
  - **Pros/Cons**: This mitigates the influence of outliers on the preprocessing step, but the PCA algorithm itself is still non-robust.
- 4.
5. **Use a Robust PCA Variant:**
- **Method**: Use a version of PCA specifically designed to be robust to outliers.
  - **Technique**: **Robust PCA** (based on Principal Component Pursuit) decomposes the data matrix D into a low-rank matrix L (the clean data) and a sparse matrix S (the outliers). It solves for L and S simultaneously.
  - **Pros/Cons**: This is the most principled approach as it handles outliers as part of the modeling process. However, it is computationally more expensive than standard PCA.
- 6.
- 

## Question 4

**How would you use PCA for data compression in a real-time streaming application?**

### Theory

Using PCA for data compression in a real-time streaming application requires an approach that can operate without storing the entire dataset and can adapt as new data arrives. This is a perfect use case for **Incremental PCA (IPCA)**. The goal is to transmit a compressed, lower-dimensional representation of the data stream, which can be reconstructed on the receiving end.

### The Framework

The system would consist of a **transmitter** (compressor) and a **receiver** (decompressor).

### **At the Transmitter (e.g., an IoT sensor):**

1. **Online Training Phase:**
  - An **Incremental PCA** model is maintained. Initially, it is trained on a buffer of the first few data batches from the stream.
  - As new data arrives in mini-batches, the `partial_fit()` method of the IPCA model is used to continuously update the principal components. This allows the model to adapt to the stream's statistics.
- 2.
3. **Real-time Compression:**
  - For each new data point or batch  $X_{\text{new}}$  from the stream:
    - The data is centered and scaled using running estimates of the mean and standard deviation.
    - It is then projected onto the currently learned top  $k$  principal components using the `transform()` method of the IPCA model. This results in a compressed vector  $Z_{\text{new}}$  of size  $k$ .
    - This much smaller vector  $Z_{\text{new}}$  is transmitted over the network.
  -
- 4.

### **At the Receiver (e.g., a central server):**

1. **Model Synchronization:** The receiver must have access to the same PCA model (the  $k$  principal components and the mean/scale vectors) as the transmitter. This model needs to be sent periodically from the transmitter to the receiver to keep it synchronized.
2. **Real-time Decompression (Reconstruction):**
  - When a compressed vector  $Z_{\text{new}}$  is received:
    - The `inverse_transform()` method of the synchronized PCA model is used to project  $Z_{\text{new}}$  back into the original high-dimensional feature space.
    - This results in a reconstructed data point  $X_{\text{reconstructed}}$ , which is an approximation of the original data.
  -
- 3.

### **Key Considerations**

- **Trade-off:** The number of components  $k$  determines the compression ratio versus the reconstruction quality. A smaller  $k$  means higher compression but more information loss.
- **Latency:** The computation at the transmitter (transform) and receiver (inverse transform) must be fast enough to keep up with the data stream's arrival rate. IPCA is designed for this.
- **Concept Drift:** If the underlying patterns in the data stream change over time, the IPCA model must be able to adapt. Some online PCA variants have "forgetting factors" to weigh recent data more heavily.

---

## Question 5

**How would you decide whether to use PCA or a classification algorithm for a given dataset?**

### Theory

This question highlights a fundamental misunderstanding of the roles of different algorithms. PCA and classification algorithms are not interchangeable; they serve completely different purposes. **PCA is an unsupervised dimensionality reduction technique**, while a **classification algorithm is a supervised learning method for prediction**.

They are not an "either/or" choice but can be used together in a pipeline.

### Correct Roles

- **PCA (Unsupervised):**
  - **Purpose:** To simplify data by reducing the number of features (dimensions). It finds the underlying structure of the features ( $X$ ) without any knowledge of a target variable ( $y$ ).
  - **Output:** A new, lower-dimensional representation of the input data.
  - **Question it Answers:** "What is the most compact way to represent the main patterns in my features?"
- 
- **Classification Algorithm (Supervised):**
  - **Purpose:** To learn a mapping from input features ( $X$ ) to a discrete target label ( $y$ ).
  - **Output:** A model that can predict the class label for new, unseen data.
  - **Question it Answers:** "Given these features, can I predict which category this sample belongs to?"
- 

### How They Work Together

The most common relationship is to use PCA as a **preprocessing step before** applying a classification algorithm.

### The Workflow:

1. You are given a dataset with features  $X$  and class labels  $y$ .
2. The goal is to build a classifier.
3. You notice that the dataset has a very high number of features, which might cause overfitting or slow training.

4. You apply **PCA** to the features X to reduce them to a smaller set of principal components X\_pca.
5. You then train your **classification algorithm** (e.g., Logistic Regression, SVM, Random Forest) using X\_pca as the input and y as the target.

### Decision Process

The decision is not "PCA or a classifier?" but rather "**Should I use PCA before my classifier?**"

You would decide to use PCA before a classifier if:

- The number of features is very large compared to the number of samples.
- The features are highly correlated (multicollinearity).
- You need to speed up the training of a computationally expensive classifier.
- You believe that reducing noise by keeping only the high-variance components will improve the model's ability to generalize.

You would decide **not** to use PCA if:

- The number of features is already small and manageable.
  - Interpretability is key, as PCA creates non-interpretable features.
  - You have reason to believe that the information needed for classification is contained in low-variance components, which PCA might discard.
- 

## Question 6

**Discuss a case where PCA helped improve model performance by reducing overfitting.**

### Theory

Overfitting occurs when a machine learning model learns the training data too well, including its noise and random fluctuations, and as a result, fails to generalize to new, unseen data. One of the primary causes of overfitting is having too many features relative to the number of training samples. PCA can help combat this by reducing the number of features, thereby simplifying the model.

### The Scenario

Let's consider a **gene expression analysis** problem.

- **Dataset:** We have a dataset of gene expression levels for thousands of genes (e.g., 20,000 features) from a small number of patients (e.g., 200 samples).
- **Goal:** To build a classifier to predict whether a patient has a certain type of cancer based on their gene expression profile.

- **The Problem (The " $p >> n$ " problem):** The number of features  $p$  (20,000) is much, much larger than the number of samples  $n$  (200). This is a classic recipe for severe overfitting. A complex model like a Support Vector Machine or a neural network could easily find a perfect separation of the training data that is based on noise, and it would perform terribly on a test set.

## How PCA Improves Performance

1. **Dimensionality Reduction:**
  - We apply PCA to the 20,000 gene expression features.
  - Genes are often co-regulated and work in pathways, meaning their expression levels are highly correlated. PCA is excellent at capturing these correlations.
  - We might find that 95% of the variance in the 20,000 features can be explained by just the first 50 principal components.
- 2.
3. **Creating a Simpler Problem:**
  - Instead of training the classifier on 20,000 noisy and correlated features, we now train it on the **50 uncorrelated principal components**.
  - The problem has been simplified from  $p=20000$  to  $p'=50$ , which is much more reasonable for a sample size of  $n=200$ .
- 4.
5. **The Result:**
  - **Reduced Overfitting:** The classifier is now forced to learn from the major patterns of variation (the signal) in the data, as the noise (contained in the discarded components) has been filtered out. This makes the model more robust and much more likely to generalize well to new patients.
  - **Improved Test Accuracy:** While the training accuracy might be lower than the "perfect" but overfitted model, the **accuracy on the test set will be significantly higher**.
- 6.

**Conclusion:** In this "fat dataset" scenario ( $p >> n$ ), PCA acts as a regularization technique. By compressing the feature space, it constrains the complexity of the problem and prevents the model from learning spurious patterns, thus improving its generalization performance.

---

## Question 7

**Give an example of how PCA might be incorrectly applied to a dataset and propose a solution.**

**Theory**

A common and critical mistake when applying PCA is **failing to standardize the data** before the analysis. This incorrect application can lead to misleading and meaningless results because PCA is sensitive to the scale of the original features.

### The Incorrect Application Scenario

- **Dataset:** A real estate dataset used to predict house prices. It contains two features:
  - `size_sqft`: The size of the house in square feet (e.g., range 1,000 - 5,000).
  - `num_bedrooms`: The number of bedrooms (e.g., range 2 - 6).
- 
- **The Goal:** Use PCA to reduce these two features to one principal component to use in a model.

### The Mistake:

A data scientist applies PCA directly to this raw, unscaled data.

1. **Variance Calculation:** The variance of `size_sqft` will be in the hundreds of thousands or millions (e.g.,  $(5000-1000)^2$  is large). The variance of `num_bedrooms` will be very small (e.g.,  $(6-2)^2$  is small).
2. **PCA Result:** PCA will find that the direction of maximum variance is almost perfectly aligned with the `size_sqft` axis. The first principal component (PC1) will be almost entirely determined by the house size. For example, PC1 might be  $0.999 * \text{size\_sqft} + 0.001 * \text{num\_bedrooms}$ .
3. **The Consequence:** The dimensionality reduction has effectively just kept the `size_sqft` feature and thrown away the information from `num_bedrooms`, regardless of how important the number of bedrooms might be for predicting the price. The result is biased and does not reflect the true combined relationship of the features.

### The Proposed Solution

The solution is simple and is a mandatory preprocessing step for PCA.

1. **Standardize the Data:**
  - Before applying PCA, use a `StandardScaler` to transform the data.
  - This will rescale both `size_sqft` and `num_bedrooms` to have a mean of 0 and a standard deviation of 1.
- 2.
3. **Apply PCA Correctly:**
  - Now, apply PCA to the **standardized** data.
  - Both features will have equal variance (1.0) and will contribute to the principal components based on their **correlation** with each other, not their arbitrary original scales.
  - The resulting first principal component will be a meaningful combination of both size and number of bedrooms, providing a much more accurate and useful summary of the house's characteristics.

4.

**Conclusion:** The cardinal rule of applying PCA is to **always standardize your data first**, unless all features are already measured in the same units and you have a specific reason to let variance dictate importance.

---

## Question 8

**Discuss how you would ensure the robustness of PCA results against variations in the dataset.**

### Theory

Ensuring the robustness of PCA results means verifying that the principal components discovered are stable and represent a genuine underlying structure in the data, rather than being an artifact of a specific sample or noise. This can be assessed using resampling techniques like bootstrapping.

### The Strategy: Bootstrapping PCA

Bootstrapping is a statistical method that involves repeatedly drawing random samples **with replacement** from the original dataset. By applying PCA to each of these bootstrap samples, we can assess the stability of the results.

#### 1. Generate Bootstrap Samples:

- Create a large number (e.g.,  $B=500$ ) of bootstrap samples from the original dataset. Each bootstrap sample will have the same size as the original dataset but will be a slightly different version due to the sampling with replacement.

2.

#### 3. Apply PCA to Each Sample:

- For each of the  $B$  bootstrap samples, perform a full PCA analysis.
- Store the results for each run, particularly the **loadings** of the principal components (the eigenvectors).

4.

#### 5. Assess Stability:

- **Component Loadings Stability:** The core of the analysis is to compare the component loadings across all bootstrap runs.
  - If the PCA results are robust, the loading vector for PC1 should be very similar across all 500 runs. The same should be true for PC2, etc.
  - We can measure this similarity by calculating the average dot product (or correlation) between the PC1 vectors from all pairs of bootstrap runs. A high average similarity indicates a very stable component.

- We can also visualize the distribution of a specific loading value (e.g., the loading of 'age' on PC1) across all runs. A narrow distribution implies stability.
  - 
  - **Eigenvalue Stability:** We can also examine the distribution of the eigenvalues for each component. A stable component will have a narrow distribution for its eigenvalue across the bootstrap samples.
- 6.

## Interpretation of Results

- **Robust Results:** If the loadings and eigenvalues for the first few principal components are highly consistent across the bootstrap samples, we can be confident that these components represent a real and stable structure in the data.
- **Non-Robust Results:** If the loadings for a component vary wildly from one bootstrap sample to the next, it suggests that this component is unstable and likely an artifact of the specific data sample. It should not be trusted or interpreted.

This bootstrapping process provides a rigorous way to validate the PCA results and gives us confidence intervals for the loadings and eigenvalues, making the analysis much more robust.

---

## Question 9

**Discuss how Randomized PCA is used and its benefits over traditional PCA.**

### Theory

**Randomized PCA** is a high-performance, approximate method for computing principal components, designed for large-scale datasets. It is based on **randomized SVD**, a probabilistic algorithm that can find a low-rank approximation of a matrix much more quickly than a traditional, exact SVD. scikit-learn uses this method when you set `svd_solver='randomized'`.

### How It Is Used

The core idea is to find a smaller subspace that captures most of the "action" of the large data matrix and then perform exact PCA on that much smaller subspace.

1. **Sketching the Subspace:** The algorithm starts by creating a "sketch" of the data. It does this by multiplying the large data matrix  $X$  by a small, random matrix  $R$ . The resulting matrix  $Y = XR$  is much smaller but preserves the essential linear structure of  $X$ .
2. **Finding an Orthonormal Basis:** It then finds an orthonormal basis  $Q$  for this sketch  $Y$ . The column space of  $Q$  serves as a very good approximation of the principal subspace of the original data.

3. **Projecting and Performing Exact PCA:** The original data  $X$  is projected onto this smaller basis  $Q$  ( $B = QTX$ ). Since  $B$  is a small matrix, a standard, exact SVD/PCA can be performed on it very quickly.
4. **Reconstructing the Components:** The final principal components for the original space are then recovered from the results of the small PCA and the basis  $Q$ .

### Benefits Over Traditional PCA

1. **Speed:** This is the primary benefit. Randomized PCA is significantly **faster** than traditional PCA (based on exact SVD), especially for high-dimensional datasets where the number of features is large but you only need to compute a small number of components. The computational complexity is much lower.
2. **Memory Efficiency:** It can be more memory-efficient as it avoids some of the large intermediate matrix computations of a full SVD.
3. **High Accuracy:** Despite being an approximation, the results are typically very close to those produced by exact PCA. For most machine learning applications, the small loss in precision is a negligible price to pay for the massive gain in speed.

### When to Use It:

Randomized PCA is the preferred method when:

- You are working with a large dataset (either many samples or many features).
  - You only need to retain a subset of the principal components ( $k$  is much smaller than the number of features).
- This is a very common scenario, making randomized PCA the default choice in many modern libraries for large-scale problems.
- 

## Question 10

**Discuss how robust PCA attempts to handle outliers and its practical implications.**

### Theory

**Robust PCA** is a modern variant of PCA specifically designed to be resilient to datasets that contain significant outliers or gross corruptions. It operates on the fundamental assumption that a data matrix  $D$  can be decomposed into the sum of two distinct matrices: a **low-rank matrix  $L$**  representing the "clean" underlying data structure, and a **sparse matrix  $S$**  representing the outliers.

$$D = L + S$$

### How It Handles Outliers

Instead of the variance-maximization objective of standard PCA, Robust PCA solves a different optimization problem, often called **Principal Component Pursuit (PCP)**.

- **Objective:** Find the matrices L and S such that the **rank of L is minimized** and the **number of non-zero elements (sparsity) of S is maximized**, subject to the constraint that  $L + S$  is close to the original data D.
- **The Intuition:**
  - By forcing L to be low-rank, the algorithm finds the main, simplified patterns in the data, which is analogous to what standard PCA does.
  - By forcing S to be sparse, the algorithm isolates the errors into a few entries. Outliers are effectively "quarantined" into the S matrix, preventing them from corrupting the estimation of the underlying structure L.
- 

## Practical Implications

1. **Simultaneous Denoising and Dimensionality Reduction:** The output of Robust PCA is twofold:
  - L: A "cleaned" version of the original data, with the outliers removed. Standard PCA can then be performed on this L matrix to get robust principal components.
  - S: A matrix that explicitly identifies the outliers. This is incredibly useful for anomaly detection.
- 2.
3. **Applications:**
  - **Background Subtraction in Video:** In a video stream, the static background is a highly correlated, low-rank phenomenon (L). Moving objects (people, cars) are transient and only affect a small part of the frame at any time, making them a sparse signal (S). Robust PCA can perfectly separate the two.
  - **Data Cleaning:** It can be used to automatically detect and separate corrupted data entries in a dataset before further analysis.
  - **Facial Recognition:** It can handle images of faces where parts are occluded (e.g., by sunglasses or a scarf). The underlying face is L, and the occlusion is S.
- 4.

**Limitation:** Robust PCA is computationally more expensive than standard PCA. However, its ability to handle grossly corrupted data makes it an invaluable tool for specific problem domains where outliers are expected.