

## 1. What are the key reasons for reducing the dimensionality of a dataset? What are the major disadvantages?

Answer:---->Disadvantages of Dimensionality Reduction It may lead to some amount of data loss. PCA tends to find linear correlations between variables, which is sometimes undesirable. PCA fails in cases where mean and covariance are not enough to define datasets.

Benefits of applying Dimensionality Reduction By reducing the dimensions of the features, the space required to store the dataset also gets reduced. Less Computation training time is required for reduced dimensions of features. Reduced dimensions of features of the dataset help in visualizing the data quickly.

## 2. What is the dimensionality curse?

Answer:---->The curse of dimensionality refers to various phenomena that arise when analyzing and organizing data in high-dimensional spaces that do not occur in low-dimensional settings such as the three-dimensional physical space of everyday experience. The expression was coined by Richard E. The curse of dimensionality, first introduced by Bellman [1], indicates that the number of samples needed to estimate an arbitrary function with a given level of accuracy grows exponentially with respect to the number of input variables (i.e., dimensionality) of the function. The large number ( $p$ ) of genes measured on a relatively small number ( $N$ ) of samples presents a difficult challenge in the analysis of DNA microarray data. This is referred to as the small  $N$ , large  $p$  problem, also called the high dimensionality problem.

## 3. Tell if its possible to reverse the process of reducing the dimensionality of a dataset? If so, how can you go about doing it? If not, what is the reason?

Answer:---->Dimensionality Reduction Techniques Feature selection. ... Feature extraction. ... Principal Component Analysis (PCA) ... Non-negative matrix factorization (NMF) ... Linear discriminant analysis (LDA) ... Generalized discriminant analysis (GDA) ... Missing Values Ratio. ... Low Variance Filter Dimensionality reduction finds a lower number of variables or removes the least important variables from the model. That will reduce the model's complexity and also remove some noise in the data. In this way, dimensionality reduction helps to mitigate overfitting. Principal Component Analysis (PCA), Factor Analysis (FA), Linear Discriminant Analysis (LDA) and Truncated Singular Value Decomposition (SVD) are examples of linear dimensionality reduction methods. Principal Component Analysis (PCA) is one of the most popular methods of dimensionality reduction as it is used for both data analysis and predictive modeling. It is a statistical method that uses an orthogonal transformation to turn observations of correlated characteristics into a set of linearly uncorrelated data.

## 4. Can PCA be utilized to reduce the dimensionality of a nonlinear dataset with a lot of variables?

Answer:---->PCA can be used to significantly reduce the dimensionality of most datasets, even if they are highly nonlinear because it can at least get rid of useless dimensions. However, if there are no useless dimensions, reducing dimensionality with PCA will lose too much information. Linear principal component analysis (PCA) can be extended to a nonlinear PCA by using artificial neural networks. But the benefit of curved components requires a careful control of the model complexity. PCA (a linear dimensionality reduction algorithm) is used to reduce this same dataset into two dimensions, the resulting values are not so well organized. Perhaps the most popular technique for dimensionality reduction in machine learning is Principal Component Analysis, or PCA for short. This is a technique that comes from the .PCA is defined as an orthogonal linear transformation that transforms the data to a new coordinate system such that the greatest variance by some scalar projection of the data comes to lie on the first coordinate (called the first principal component), the second greatest variance on the second coordinate, and so on. field of linear algebra and can be used as a data preparation technique to create a projection of a dataset prior to fitting a model

5. Assume you're running PCA on a 1,000-dimensional dataset with a 95 percent explained variance ratio. What is the number of dimensions that the resulting dataset would have?

Answer:---->If you are running PCA on a 1,000-dimensional dataset with a 95% explained variance ratio, the resulting dataset would have a number of dimensions equal to the number of components that are needed to capture 95% of the variance in the original dataset. This number can be computed using the explained variance ratio of each component, which can be obtained by running PCA on the original dataset.

For example, if the first component has an explained variance ratio of 50%, the second component has an explained variance ratio of 25%, and the third component has an explained variance ratio of 20%, then the first three components would capture 95% of the variance in the original dataset ( $50\% + 25\% + 20\% = 95\%$ ). Therefore, the resulting dataset would have 3 dimensions.

It is important to note that the number of dimensions in the resulting dataset will depend on the specific explained variance ratios of the components, and may not necessarily be equal to the number of dimensions needed to capture 95% of the variance in all cases.

6. Will you use vanilla PCA, incremental PCA, randomized PCA, or kernel PCA in which situations?

Answer:---->There are several variations of principal component analysis (PCA) that can be used for different purposes. Here is a brief summary of when each variation may be appropriate:

**Vanilla PCA:** This is the standard form of PCA, which computes the principal components of a dataset by performing singular value decomposition (SVD) on the data matrix. It is suitable for situations where the dataset is small enough to fit in memory, and the computation time is not a concern.

**Incremental PCA:** This variation of PCA is designed to handle large datasets that cannot be fit in memory. It performs PCA incrementally by processing small batches of the data at a time, and can be used when the data is too large to fit in memory all at once.

**Randomized PCA:** This variation of PCA uses randomized algorithms to approximate the principal components of a dataset. It is faster than vanilla PCA, but may not be as accurate. It is suitable for situations where computation time is a concern, and a good approximation of the principal components is sufficient.

**Kernel PCA:** This variation of PCA is used to perform nonlinear dimensionality reduction by projecting the data onto a higher-dimensional feature space using a kernel function. It is suitable for situations where the relationships between the features in the original dataset are nonlinear, and linear PCA is not sufficient.

In general, the choice of which variation of PCA to use will depend on the specific characteristics of the dataset and the computational resources available.

7. How do you assess a dimensionality reduction algorithm's success on your dataset?

Answer:---->There are several ways to assess the success of a dimensionality reduction algorithm on a dataset. Some common measures include:

**Reconstruction error:** The reconstruction error measures the difference between the original dataset and the reconstructed dataset after dimensionality reduction. A lower reconstruction error indicates that the dimensionality reduction algorithm was able to preserve more of the important information in the original dataset.

**Explained variance ratio:** The explained variance ratio measures the amount of variance in the original dataset that is captured by the reduced dataset. A higher explained variance ratio indicates that the dimensionality reduction algorithm was able to capture more of the important information in the original dataset.

**Visualization:** Visualizing the data after dimensionality reduction can provide insight into how well the algorithm was able to preserve the structure of the original dataset. For example, if the data points in the reduced dataset are well-separated and distinct, it may indicate that the algorithm was successful in preserving the structure of the original data.

Classification or regression performance: If the goal of the dimensionality reduction is to improve the performance of a supervised learning task (such as classification or regression), then the performance of the learning algorithm on the reduced dataset can be used as a measure of the success of the dimensionality reduction algorithm. A higher performance on the reduced dataset compared to the original dataset may indicate that the dimensionality reduction algorithm was successful in removing noise and improving the signal-to-noise ratio of the data.

It is important to note that different measures may be more or less relevant depending on the specific characteristics of the dataset and the goals of the dimensionality reduction.

## ▼ 8. Is it logical to use two different dimensionality reduction algorithms in a chain?

Answer:---->It is generally possible to use two different dimensionality reduction algorithms in a chain, where the output of one algorithm is used as the input to the next. This approach can sometimes be useful if the first algorithm is not able to effectively reduce the dimensionality of the dataset, and a second algorithm is needed to further reduce the dimensionality.

However, it is important to keep in mind that using multiple dimensionality reduction algorithms in a chain can increase the risk of losing important information in the dataset. Each algorithm has its own set of assumptions and constraints, and applying multiple algorithms in a chain may result in a dataset that is not representative of the original data.

Therefore, it is generally advisable to carefully consider the goals and limitations of the dimensionality reduction task, and choose the most appropriate algorithm or combination of algorithms based on these considerations. In some cases, using multiple dimensionality reduction algorithms in a chain may be beneficial, but in other cases it may be more appropriate to use a single algorithm or a different combination of algorithms.

[Colab paid products](#) - [Cancel contracts here](#)

