

Reducción de dimensionalidad

Dr. Mauricio Toledo-Acosta

Diplomado Ciencia de Datos con Python

Table of Contents

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

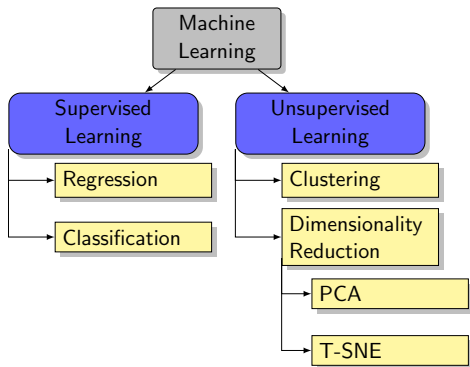
Introducción

PCA

1 Introducción

2 PCA

Introducción



Reducción de dimensionalidad

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

Introducción

PCA

Reducción de dimensionalidad

La reducción de la dimensionalidad es la transformación de los datos de un espacio de alta dimensión a un espacio de baja dimensión, de manera que la representación de baja dimensión conserve algunas propiedades significativas de los datos originales.

Reducción de dimensionalidad

Reducción de dimensionalidad

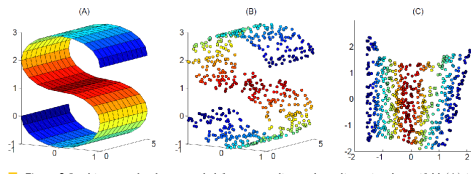
Reducción de dimensionalidad

Introducción

PCA

Reducción de dimensionalidad

La reducción de la dimensionalidad es la transformación de los datos de un espacio de alta dimensión a un espacio de baja dimensión, de manera que la representación de baja dimensión conserve algunas propiedades significativas de los datos originales.



Reducción de dimensionalidad

Reducción de dimensionalidad

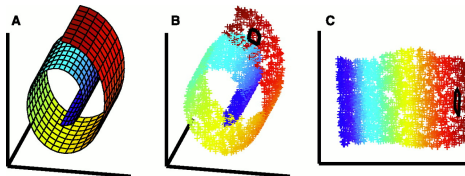
Reducción de dimensionalidad

Introducción

PCA

Reducción de dimensionalidad

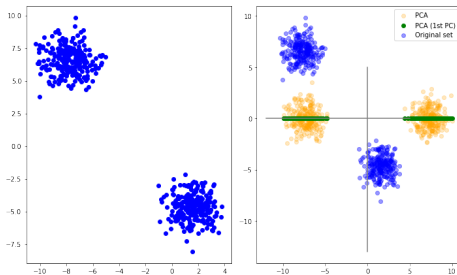
La reducción de la dimensionalidad es la transformación de los datos de un espacio de alta dimensión a un espacio de baja dimensión, de manera que la representación de baja dimensión conserve algunas propiedades significativas de los datos originales.



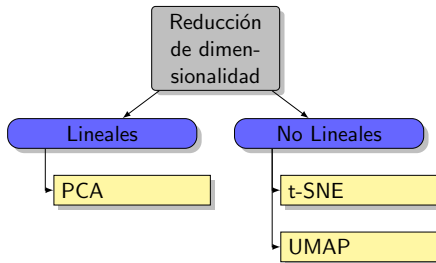
Reducción de dimensionalidad

Reducción de dimensionalidad

La reducción de la dimensionalidad es la transformación de los datos de un espacio de alta dimensión a un espacio de baja dimensión, de manera que la representación de baja dimensión conserve algunas propiedades significativas de los datos originales.



◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ 🔍 ↺ ↻



Utilidad

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

Introducción

PCA

¿Para qué sirve la reducción de dimensionalidad?

- Visualización.

Utilidad

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

Introducción

PCA

¿Para qué sirve la reducción de dimensionalidad?

- Visualización.
- Extraer la información más importante de los datos.

Utilidad

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

Introducción

PCA

¿Para qué sirve la reducción de dimensionalidad?

- Visualización.
- Extraer la información más importante de los datos.
- Obtener features para fines de clasificación.

Table of Contents

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

Introducción

PCA

1 Introducción

2 PCA

PCA

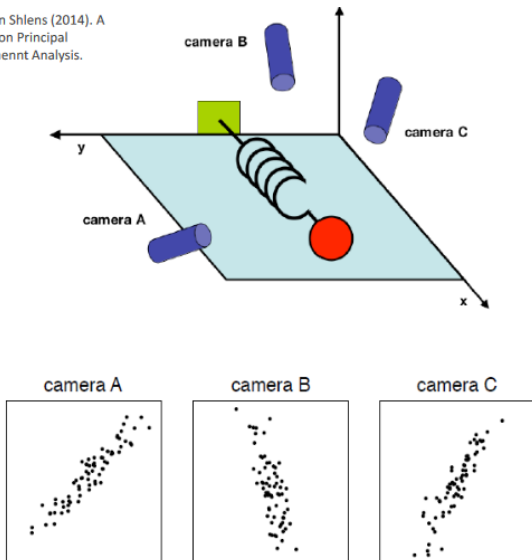
Reducción de dimensionalidad

Reducción de dimensionalidad

Introducción

PCA

Jonathon Shlens (2014). A tutorial on Principal Component Analysis.



PCA

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

Introducción

PCA

PCA puede pensarse como el **ajuste de un elipsoide D -dimensional al conjunto de datos**, donde cada eje del elipsoide representa una componente principal. Si algún eje del elipsoide es pequeño, entonces la varianza a lo largo de ese eje también es pequeña.

PCA

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

Introducción

PCA

PCA puede pensarse como el **ajuste de un elipsoide D -dimensional al conjunto de datos**, donde cada eje del elipsoide representa una componente principal. Si algún eje del elipsoide es pequeño, entonces la varianza a lo largo de ese eje también es pequeña.

PCA transforma los datos a un nuevo sistema de coordenadas de tal manera que la mayor varianza se sitúa en la primera coordenada, la segunda mayor varianza en la segunda coordenada, y así sucesivamente.

PCA

Reducción de dimensionalidad

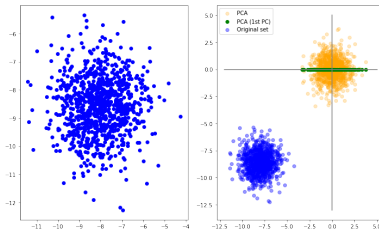
Reducción de dimensionalidad

Introducción

PCA

PCA puede pensarse como el **ajuste de un elipsoide D -dimensional al conjunto de datos**, donde cada eje del elipsoide representa una componente principal. Si algún eje del elipsoide es pequeño, entonces la varianza a lo largo de ese eje también es pequeña.

PCA transforma los datos a un nuevo sistema de coordenadas de tal manera que la mayor varianza se sitúa en la primera coordenada, la segunda mayor varianza en la segunda coordenada, y así sucesivamente.

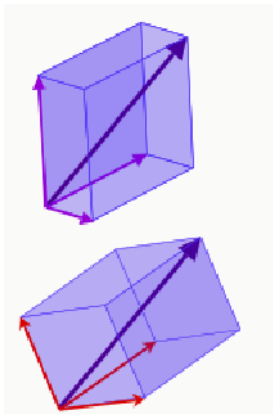


Cambios de Base

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

Introducción
PCA



- Las coordenadas de un punto en son los coeficientes de los vectores canónicos unitarios

$$e_1 = (1, 0, \dots, 0),$$

...

$$e_D = (0, 0, \dots, 1).$$

- Todo punto puede ser expresado en una infinidad de bases.
- Para cambiar de base hay que multiplicar el vector por la matriz

$$\begin{pmatrix} v_1^{(1)} & \dots & v_1^{(D)} \\ \vdots & \ddots & \vdots \\ v_D^{(1)} & \dots & v_D^{(D)} \end{pmatrix}$$

PCA

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

Introducción

PCA

¿Cómo obtenemos la nueva base de coordenadas para PCA? Esta base de vectores deben ser las direcciones de máxima varianza, es decir, las **componentes principales**.

PCA

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

Introducción

PCA

¿Cómo obtenemos la nueva base de coordenadas para PCA? Esta base de vectores deben ser las direcciones de máxima varianza, es decir, las **componentes principales**.

Consideremos la matriz de covarianza

$$C_X = \frac{1}{N} X \cdot X^T.$$

PCA

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

Introducción

PCA

¿Cómo obtenemos la nueva base de coordenadas para PCA? Esta base de vectores deben ser las direcciones de máxima varianza, es decir, las **componentes principales**.

Consideremos la matriz de covarianza

$$C_X = \frac{1}{N} X \cdot X^T.$$

- Los términos altos en la diagonal corresponden a una varianza alta.
- Los términos altos fuera de la diagonal corresponden a una redundancia alta.

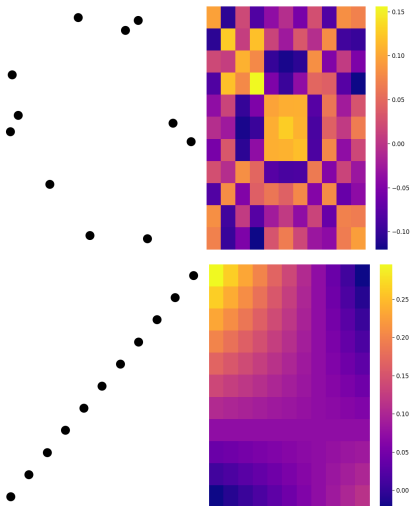
Matriz de covarianza

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

Introducción

PCA



Matriz de Covarianza

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

Introducción

PCA

De acuerdo con lo anterior, para nuestra matriz de covarianza C_X deseamos:

- Minimizar la redundancia, medida por la magnitud de la covarianza.
- Maximizar la varianza.

Matriz de Covarianza

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

Introducción

PCA

De acuerdo con lo anterior, para nuestra matriz de covarianza C_X deseamos:

- Minimizar la redundancia, medida por la magnitud de la covarianza.
- Maximizar la varianza.

Esto, en términos de matrices, quiere decir **Diagonalizar**. Es decir, encontrar matrices P y D tales que

$$D = P \cdot C_X \cdot P^T$$

Matriz de Covarianza

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

Introducción

PCA

De acuerdo con lo anterior, para nuestra matriz de covarianza C_X deseamos:

- Minimizar la redundancia, medida por la magnitud de la covarianza.
- Maximizar la varianza.

Esto, en términos de matrices, quiere decir **Diagonalizar**. Es decir, encontrar matrices P y D tales que

$$D = P \cdot C_X \cdot P^T$$

Esto se hace encontrando los eigenvector y eigenvalores de C_X .

Eigenvalores y Eigenvectores

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

Introducción

PCA

Cada matriz M de tamaño $n \times n$ se puede ver como una transformación del espacio \mathbb{R}^n en el espacio \mathbb{R}^n . Es decir, toma un punto $p \in \mathbb{R}^n$ y devuelve otro vector $M \cdot p \in \mathbb{R}^n$. ¿Cómo es este punto respecto al inicial?

Eigenvalores y Eigenvectores

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

Introducción

PCA

Cada matriz M de tamaño $n \times n$ se puede ver como una transformación del espacio \mathbb{R}^n en el espacio \mathbb{R}^n . Es decir, toma un punto $p \in \mathbb{R}^n$ y devuelve otro vector $M \cdot p \in \mathbb{R}^n$. ¿Cómo es este punto respecto al inicial?

Hay vectores especiales, dependiendo de M , que son transformados en un múltiplo de ellos mismos.

$$A = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}$$

$$A \cdot \begin{pmatrix} 1 \\ -1 \end{pmatrix} = \begin{pmatrix} 1 \\ -1 \end{pmatrix} = 1 \cdot \begin{pmatrix} 1 \\ -1 \end{pmatrix}$$

$$A \cdot \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 3 \\ 3 \end{pmatrix} = 3 \cdot \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

Eigenvalores y Eigenvectores

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

Introducción

PCA

Cada matriz M de tamaño $n \times n$ se puede ver como una transformación del espacio \mathbb{R}^n en el espacio \mathbb{R}^n . Es decir, toma un punto $p \in \mathbb{R}^n$ y devuelve otro vector $M \cdot p \in \mathbb{R}^n$. ¿Cómo es este punto respecto al inicial?

Hay vectores especiales, dependiendo de M , que son transformados en un múltiplo de ellos mismos.

$$A = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}$$
$$A \cdot \begin{pmatrix} 1 \\ -1 \end{pmatrix} = \begin{pmatrix} 1 \\ -1 \end{pmatrix} = 1 \cdot \begin{pmatrix} 1 \\ -1 \end{pmatrix}$$
$$A \cdot \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 3 \\ 3 \end{pmatrix} = 3 \cdot \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

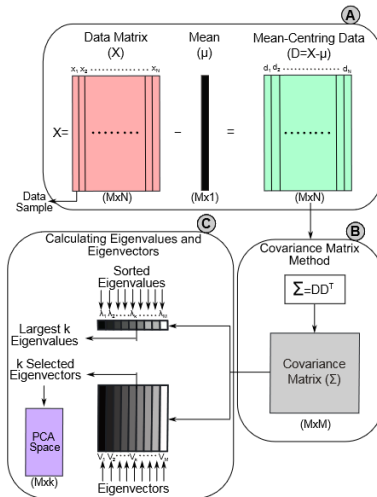
El proceso

Reducción de dimensionalidad

Reducción de dimensionalidad

Introducción

PCA



Conexiones entre la geométricas y la estadística

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

Introducción

PCA

- Promedio/Centroide = Media.
- Norma = Varianza.
- Ángulo entre vectores = Covarianza.

Ventajas y Desventajas

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

Introducción

PCA

Ventajas

- Permite eliminar variables correlacionadas.
- Permite la visualización de datos.
- Puede ayudar a reducir el overfitting.

Ventajas y Desventajas

Reducción de
dimensionali-
dad

Reducción de
dimensionali-
dad

Introducción

PCA

Ventajas

- Permite eliminar variables correlacionadas.
- Permite la visualización de datos.
- Puede ayudar a reducir el overfitting.

Desventajas

- Suele requiere de un escalamiento de datos antes.
- Se pierde información.
- Perdemos la interpretabilidad de las variables de entrada.