



Detalles de la derivación e implementación de algoritmo básico para la reducción de dimensionalidad con PCA

Details on the algorithm for dimensionality reduction via Principal Components Analysis

Detalhes da derivação e implementação do algoritmo básico para redução de dimensionalidade com PCA

Zenaida Natividad Castillo-Marrero ^I
zenaida.castillo@esPOCH.edu.ec
<https://orcid.org/0000-0002-4424-8652>

Gustavo Adolfo Colmenares-Pacheco ^{II}
gcolmenares@yachaytech.edu.ec
<https://orcid.org/0000-0003-4789-0859>

Ramón Antonio Abancin-Ospina ^{III}
ramon.abancin@esPOCH.edu.ec
<https://orcid.org/0000-0002-2417-6671>

Víctor Oswaldo Cevallos-Vique ^{IV}
victor.cevallos@esPOCH.edu.ec
<https://orcid.org/0000-0001-5525-5818>

Correspondencia: zenaida.castillo@esPOCH.edu.ec

Ciencias Técnicas y Aplicadas
Artículo de Investigación

***Recibido:** 04 de enero de 2022 ***Aceptado:** 31 de enero de 2022 * **Publicado:** 21 de febrero de 2022

- I. Escuela Superior Politécnica de Chimborazo (ESPOCH), Facultad de Ciencias, Escuela de Matemática, Grupo CITED, Riobamba, Ecuador.
- II. Universidad de Investigación de Tecnología Experimental Yachay, Escuela de Ciencias Matemáticas y Computacionales, Urcuquí, 060150, Ecuador.
- III. Escuela Superior Politécnica de Chimborazo (ESPOCH), Facultad de Ciencias, Escuela de Matemática, Grupo CITED, Riobamba, Ecuador.
- IV. Escuela Superior Politécnica de Chimborazo (ESPOCH), Facultad de Administración y Empresas, Escuela de Finanzas, Grupo CITED, Riobamba, Ecuador.

Resumen

Los problemas que requieren de análisis de datos, a menudo son difíciles de resolver, debido principalmente a la cantidad de variables involucradas en el modelo matemático. Los científicos de datos generalmente trabajan con millones de variables para hacer estimaciones que soportan decisiones importantes. En el procesamiento digital de imágenes, por ejemplo, el número de puntos que representan píxeles en tres dimensiones podría muy grande en imágenes a color. En estos casos, el costo computacional que requiere el manejo de estos datos puede resultar inaceptable, y la reducción de dimensionalidad de estos datos se hace necesaria. Aún cuando se cuente con la tecnología adecuada, la reducción en tiempo de cómputo siempre es deseable. El manejo de datos con alta dimensionalidad, el análisis y la interpretación se dificulta y en el caso de imágenes, su visualización podría verse afectada considerando las limitaciones de memoria. En la mayoría de los casos, estos datos son redundantes, y la información importante puede revelarse con solo parte de los mismos. La reducción de dimensionalidad es el proceso mediante el cual se descarta parte de la data que no aporta información relevante; y uno de los métodos más usados en todos los ámbitos es el análisis de componentes principales, o PCA, el cual se basa en el cálculo de algunos autovalores de la matriz de covarianzas de los datos. En este trabajo revisamos las herramientas matemáticas detrás del análisis del PCA, y detallamos los pasos de un algoritmo para reducir dimensionalidad que luego es implementado en Matlab. Se presenta también un ejemplo de aplicación para ilustrar el proceso.

Palabras claves: PCA; Reducción de dimensionalidad; Autovalores; Análisis de datos.

Abstract

In data analysis, we often deal with millions of variables to make estimations for decisions. Also, in digital image processing the number of data points representing pixels in three dimensions could become very large for color images. In all these cases computational costs for solving associated problems could not be affordable, and a reduction of dimensionality it is necessary. When working with high dimensional data, analysis and interpretation of results is difficult, and visualization of images could be prohibited in terms of memory. Frequently the data contains redundancy, and the important information can be found just with part of the data. In a color image, for example, with four channels, red, blue, green and grey-scale channel, one could combine the first three channels to get the information of the gray-scale channel. Dimensional

reduction is a technique that allow us to work with a reduced part of the data without lost losing important information, and one of the methods used for dimensional reduction is the principal analysis component or PCA, which is based in the computation of a few eigenvalues of the covariance matrix. In this work we review the math tools and the steps to generate a classical algorithm for dimensionality reduction using PCA. We present mathematical elements and concepts involved in a principal component analysis algorithm for the projection of high dimensional data in a lower dimensional subspace. A practical example is presented to show how these mathematical tools work together to produce similar results in lower dimensional spaces.

Keywords: PCA; Dimensionality reduction; Eigenvalues; Data analysis.

Resumo

Problemas que requerem análise de dados são muitas vezes difíceis de resolver, principalmente devido ao número de variáveis envolvidas no modelo matemático. Os cientistas de dados normalmente trabalham com milhões de variáveis para fazer estimativas que apoiam decisões importantes. No processamento digital de imagens, por exemplo, o número de pontos representando pixels em três dimensões pode ser muito grande em imagens coloridas. Nesses casos, o custo computacional necessário para lidar com esses dados pode ser inaceitável, e a redução da dimensionalidade desses dados torna-se necessária. Mesmo quando a tecnologia apropriada está disponível, a redução do tempo de computação é sempre desejável. O manuseio de dados com alta dimensionalidade, análise e interpretação é difícil e, no caso de imagens, sua visualização pode ser afetada por limitações de memória. Na maioria dos casos, esses dados são redundantes e informações importantes podem ser reveladas apenas com parte deles. A redução de dimensionalidade é o processo pelo qual parte dos dados que não fornecem informações relevantes é descartada; e um dos métodos mais utilizados em todas as áreas é a análise de componentes principais, ou PCA, que se baseia no cálculo de alguns autovalores da matriz de covariância dos dados. Neste artigo, revisamos as ferramentas matemáticas por trás da análise de PCA e detalhamos as etapas de um algoritmo para reduzir a dimensionalidade que é implementado no Matlab. Um exemplo de aplicação também é apresentado para ilustrar o processo.

Palavras-chave: PCA; Redução de dimensionalidade; Autovalores; Analise de dados.

Introducción

La reducción de dimensionalidad no es más que la determinación de un número de variables mínimo con el cual se puedan representar los datos, y consiste en eliminar variables que no sean relevantes o que no aporten suficiente información. Actualmente existen técnicas o algoritmos para hacer esta tarea, la cual ha pasado a ser una actividad obligada en disciplinas como la ciencia de datos, una vez que simplifica el manejo de los datos y al mismo tiempo reduce el costo computacional en tiempo de procesamiento y utilización de memoria. Una técnica usada para reducir la dimensionalidad es eficiente mientras logra esta reducción de costos al mismo tiempo que mantiene una buena representación del modelo.

En la actualidad muchas aplicaciones dependen del manejo apropiado de datos, y a menudo los datos suelen presentarse en muchas dimensiones, dependiendo de la cantidad de variables involucradas, que en principio se justifican para la toma de decisiones; por lo tanto se espera que al tener representación en muchas dimensiones la información sea más completa. Sin embargo, el manejo y la interpretación de los datos se hace más complicado, sin contar el costo computacional que a veces es imposible de afrontar, véase por ejemplo (James, 2013)

La reducción de dimensionalidad explota la estructura del problema respetando la correlación entre variables, lo cual permite un manejo más eficiente de los datos sin pérdida de la información relevante.

Este trabajo está orientado a describir en forma simple para lectores que se inician en el tópico, las herramientas matemáticas que se usan en el proceso de la reducción lineal de dimensionalidad a través del análisis de componentes principales o PCA (por sus siglas en inglés). Se mostrará la derivación y posterior implementación en Matlab del algoritmo clásico que define esta técnica para el análisis de datos.

En las próximas secciones estaremos definiendo conceptos básicos asociados al análisis de datos a través de la estadística, con el objetivo de identificar el aporte del análisis de componentes principales al manejo eficiente y significativo de datos con altas dimensiones.

Medidas estadísticas empíricas de los datos

Generalmente, el manejo de los datos, a través de una muestra de ellos, a fin de obtener información de interés, comienza por establecer parámetros estadísticos como la media y la varianza, que nos dan información sobre la dispersión de los datos. En esta sección se presentan

estos conceptos que miden la variabilidad de los datos. Comenzaremos considerando datos en una dimensión (1D) y luego extenderemos los conceptos a dos 2D o más dimensiones.

Sea $D = \{x_1, x_2, \dots, x_N\}$, un conjunto de datos u observaciones en una dimensión donde cada $x_i \in \mathbb{R}$, se define la media del conjunto como:

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

En este caso, cada x_i representa una característica o variable que se desea estudiar.

La varianza se define como:

$$V[D] = s^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$$

En algunas aplicaciones se divide por $N - 1$ para que sea insesgado.

La varianza, como estadístico de dispersión, mide la variabilidad de los datos de su medida central o media aritmética. Una medida que suele darnos más información sobre los datos es la desviación estándar s , o la raíz cuadrada de la varianza, la cual indica el rango de la variabilidad en los datos. La desviación estándar se indica en las mismas unidades de la media, mientras que la varianza se especifica en unidades al cuadrado por lo que no es tan fácil conjeturar sobre la información que brinda.

Se entiende por población a un conjunto bien definido de elementos; por ejemplo, los habitantes de una comunidad o los estudiantes de una Universidad. Si por cada elemento de una población se mide un conjunto de variables estaremos hablando de estadística multivariante, y los estadísticos como la media y la varianza se extienden apropiadamente. Por ejemplo, en el caso de los estudiantes de una Universidad pudiéramos recoger datos como la edad, género, promedio de calificaciones, carrera que cursa, etc., y por cada estudiante obtener esta información con el fin de analizarla para tomar alguna decisión gerencial.

Cada característica o variable puede tomar valores cuantitativos o cualitativos; sin embargo, en la práctica los variables cualitativas como color o género podrán representarse en una escala numérica; de esta manera si tenemos n individuos y medimos k características o variables podemos utilizar una matriz para representar estos datos:

$$X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1k} \\ x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix} = \begin{bmatrix} \vec{x}_1^T \\ \vec{x}_2^T \\ \vdots \\ \vec{x}_n^T \end{bmatrix}$$

Donde \vec{x}_i^T es el vector que representa las características del i -ésimo individuo. En este caso tendremos un vector de medias con k componentes, cada una representando la media de los valores de la variable respectiva.

$$\vec{\bar{x}} = \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_k \end{bmatrix}, \text{ donde } \bar{x}_i = \frac{1}{N} \vec{x}_i^T \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

La dispersión o variabilidad de los datos en el caso multivariable puede ser analizada a través de las relaciones lineales entre las variables y representada mediante una matriz llamada matriz de varianzas y covarianzas, o simplemente matriz de covarianza.

$$S = \begin{bmatrix} s_1^2 & s_{12} & \cdots & s_{1k} \\ s_{21} & s_2^2 & \cdots & s_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ s_{k1} & s_{k2} & \cdots & s_{kk} \end{bmatrix}$$

Las entradas de la diagonal de esta matriz son las varianzas de cada variable, mientras que las entradas fuera de la diagonal son las covarianzas o los indicadores bivariantes. La entrada S_{ij} ($i \neq j$) de esta matriz mide la variación conjunta de las variables x_i , y x_j ; si S_{ij} es positiva significa que a valores altos/bajos de la variable x_i podemos esperar valores altos/bajos de x_j . Si el valor de S_{ij} es negativo significa que la covariación es en sentido inverso. Si la covarianza es cero no se puede establecer una relación entre la variabilidad e las variables.

El cálculo de S también puede realizarse como:

$$S = \frac{1}{N} \sum_{i=1}^N (X_i - \bar{x})(X_i - \bar{x})^T \in \mathbb{R}^{D \times D}$$

Para efectos prácticos es importante reconocer que la matriz S es simétrica y semidefinida positiva, lo cual garantiza, entre otras cosas, que:

- i. Para cualquier vector v de \mathbb{R}^D : $v^T S v \geq 0$.
- ii. Todos los autovalores de S son positivos o cero.

La covarianza de dos variables x , y , con medias \bar{x} , \bar{y} respectivamente, puede calcularse directamente como:

$$S = Cov(x,y) = \frac{\sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{N}$$

Cuando los datos están centrados alrededor de la media, la esperanza o valor esperado de la variable aleatoria es $E[x]=0$. En este caso la matriz de covarianzas es:

$$S = \frac{1}{N} \sum_{i=1}^N X_i X_i^T \in \mathbb{R}^{D \times D}$$

Otro estadístico importante que señala en qué proporción una variable y está explicada por la influencia lineal de otra variable x , es el coeficiente de correlación lineal, denotado por $r(x,y)$.

$$r(x,y) = \frac{Cov(x,y)}{s_x s_y}$$

El coeficiente de correlación lineal toma valores en $[-1,1]$, si $r = 1$ se dice que las variables están directamente relacionadas si una crece/decrece la otra crece/decrece; y si $r = -1$, la relación es contraria, si una crece/decrece la otra decrece/crece. Cuando $r = 0$ decimos que las variables no están correlacionadas.

Una explicación detallada de estos estadísticos con ejemplos incluidos puede ser hallada en (Salazar, 2017, Rencher, 1994).

Transformaciones lineales de los datos

En muchas aplicaciones se desea conocer lo que pasa con estos estadísticos si a los datos se les suma una cantidad o si se multiplican por un escalar; es decir, si son sometidos a transformaciones de desplazamiento o de compresión/estiramiento. Para analizar matemáticamente esta situación, consideremos un conjunto de datos $\mathbf{D} = \{x_1, x_2, \dots, x_N\}$, donde

$$x_i \in \mathbb{R} \text{ (1-dimensión), } \mathbf{E}[\mathbf{D}] = \mu, \text{ y } \mathbf{V}[\mathbf{D}] = \sigma^2.$$

Si cada x_i se modifica como $x'_i = \alpha x_i + c$, con α y c dados, entonces,

$$\mathbf{E}[\alpha \mathbf{D}] = \alpha \mathbf{E}[\mathbf{D}] \wedge \mathbf{E}[\mathbf{D} + c] = \mathbf{E}[\mathbf{D}] + c \Rightarrow \mathbf{E}[\alpha \mathbf{D} + c] = \alpha \mathbf{E}[\mathbf{D}] + c$$

Ilustremos este resultado con un ejemplo, sea $\mathbf{D} = \{1, 2, 3\}$, y tomemos $\alpha = 3$ y $c = 2$.

Entonces

$$\text{i) } \mathbf{E}[\mathbf{D}] = (1 + 2 + 3)/3 = 2 \text{ y } \mathbf{E}[\alpha \mathbf{D}] = \mathbf{E}[3\mathbf{D}] = (3 + 6 + 9)/3 = 6 = 3\mathbf{E}[\mathbf{D}].$$

$$\text{ii) } \mathbf{E}[\mathbf{D} + c] = \mathbf{E}[\mathbf{D} + 2] = (3 + 4 + 5)/3 = 4 = \mathbf{E}[\mathbf{D}] + 2.$$

$$\text{iii) } \mathbf{E}[\alpha \mathbf{D} + c] = \mathbf{E}[3\mathbf{D} + 2] = (5 + 8 + 11)/3 = 8 = 3\mathbf{E}[\mathbf{D}] + 2.$$

Una situación similar se presenta con la varianza.

$$V[\alpha D] = \alpha^2 V[D] \wedge V[D + c] = V[D] \Rightarrow V[\alpha D + c] = \alpha^2 V[D]$$

Ejemplificamos estas propiedades con los datos del ejemplo anterior.

$$\text{i) } V[D] = ((-1)^2 + 0^2 + 1^2)/3 = 2/3, \quad V[\alpha D] = V[3D] = (9 + 0 + 9)/3 = 3^2 V[D].$$

$$\text{ii) } V[D + c] = V[D + 2] = ((-1)^2 + 0^2 + 1^2)/3 = 2/3 = V[D].$$

$$\text{iii) } V[\alpha D + c] = V[3D + 2] = ((-3)^2 + 0^2 + 3^2)/3 = 6 = 3^2 V[D].$$

El detalle de estas propiedades y su deducción puede ser hallado en (De la Puente, 2018, Rencher 1994).

Algunos conceptos de espacios vectoriales

A continuación revisamos brevemente conceptos del álgebra lineal, dentro del tópico de espacios vectoriales, cuyo entendimiento nos permitirá entender cómo proyectar en espacios de menor dimensión.

Producto escalar: Sean $\mathbf{x} = (x_1, x_2, \dots, x_n)$ y $\mathbf{y} = (y_1, y_2, \dots, y_n)$ dos vectores de \mathbb{R}^n , se define su producto escalar o producto punto como:

$$\mathbf{x}^T \mathbf{y} = \sum_{i=1}^n x_i y_i.$$

El producto escalar es un producto interior, o producto interno, este último definido como una función que toma dos elementos de un espacio vectorial V y devuelve un número real, usualmente se denota como $\langle \cdot, \cdot \rangle: V \times V \rightarrow \mathbb{R}$, y satisface las siguientes propiedades:

- I. Simetría: para todo par de vectores \mathbf{x} , \mathbf{y} se cumple que $\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle$.
- II. Definida positiva: Para $\mathbf{x} \neq \mathbf{0}$ se cumple que $\langle \mathbf{x}, \mathbf{x} \rangle = 0$, y $\langle \mathbf{0}, \mathbf{0} \rangle = 0$.
- III. Bilineal: Para todo $\mathbf{x}, \mathbf{y}, \mathbf{z} \in V$, y $\lambda \in \mathbb{R}$ se cumple que:

$$\circ \langle \lambda \mathbf{x} + \mathbf{y}, \mathbf{z} \rangle = \lambda \langle \mathbf{x}, \mathbf{z} \rangle + \langle \mathbf{y}, \mathbf{z} \rangle$$

$$\circ \langle \mathbf{x}, \lambda \mathbf{y} + \mathbf{z} \rangle = \lambda \langle \mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{x}, \mathbf{z} \rangle$$

La norma de un vector \mathbf{x} está vinculada al producto interno $\|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$, y podría ser diferente para dos funciones que definan productos internos distintos. Por ejemplo, en \mathbb{R}^2 tomando dos vectores genéricos $\mathbf{x} = (x_1, x_2)$, $\mathbf{y} = (y_1, y_2)$ pudiéramos definir el siguiente producto interno:

$$\langle \mathbf{x}, \mathbf{y} \rangle = x_1 y_1 - \frac{1}{2} x_1 y_2 - \frac{1}{2} x_2 y_1 + x_2 y_2$$

Una inspección de esta función nos hará concluir que en efecto satisface las propiedades I, II y III, mencionadas previamente, que definen un producto interno. En este producto interno, la longitud del vector $\mathbf{x} = (1, 1)$ será $\|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle} = 1$, mientras que en el producto punto

$$\|\mathbf{x}\| = \sqrt{\mathbf{x}^T \mathbf{x}} = 2.$$

Otro ejemplo de funciones que definen un producto interno es la covarianza, que por definición es bilineal, simétrica y definida positiva.

De igual manera podemos predecir que la distancia entre vectores depende del producto interno que se utilice. Si usamos el producto punto, más generalmente conocido, obtenemos lo que se conoce como la distancia euclídea.

En general, dado un producto interno $\langle \cdot, \cdot \rangle$ la distancia entre dos vectores \mathbf{x}, \mathbf{y} se define como:

$$\text{dist}(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\| = \sqrt{\langle \mathbf{x} - \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle}$$

Otro aspecto geométrico a considerar en el análisis de datos es la orientación de los datos, la cual como veremos en la próxima sección puede establecerse mediante direcciones dadas por vectores. Así, la medida del ángulo entre dos vectores puede decirnos qué tan similar son sus direcciones, y esta medida también está relacionada con el producto interno y la norma de los vectores; por ejemplo:

$$\cos(\theta) = \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\|\mathbf{x}\| \|\mathbf{y}\|}$$

Ortogonalidad: Se dice que dos vectores, distintos de cero, son ortogonales si su producto interno es cero. Por ejemplo, los vectores $\mathbf{x} = (1, 2)$ y $\mathbf{y} = (-2, 1)$ son ortogonales ya que su producto escalar es cero; haciendo notar que en el producto interno definido en (1) estos vectores no serían ortogonales ya que

$$\langle \mathbf{x}, \mathbf{y} \rangle = x_1 y_1 - \frac{1}{2} x_1 y_2 - \frac{1}{2} x_2 y_1 + x_2 y_2 = -2 - 0.5 + 2 + 2 = 1.5.$$

Una base de un espacio vectorial \mathbf{V} es un conjunto de vectores de \mathbf{V} que son linealmente independientes y generan a cualquier otro vector de \mathbf{V} a través de alguna combinación lineal. Por ejemplo, en \mathbb{R}^2 los conjuntos \mathbf{B}_1 y \mathbf{B}_2 cumplen con tales requisitos.

$$\mathbf{B}_1 = \{(1, 0), (0, 1)\}, \mathbf{B}_2 = \{(1, 2), (2, 1)\}$$

La base \mathbf{B}_1 del ejemplo, conocida como la base canónica de \mathbb{R}^2 , contiene dos vectores ortogonales entre sí y de norma 1, por lo que decimos que es una base ortonormal de \mathbb{R}^2 . Existen otras bases con estas características que resultan provechosas cuando se hacen proyecciones a

espacios de menor dimensión. Todas las bases de un espacio vectorial tienen el mismo número de vectores, y a este número se le conoce como la dimensión del espacio. El concepto se puede extender a \mathbb{R}^n y también a espacios vectoriales de funciones.

Una base ortonormal de vectores en \mathbb{R}^n es un subconjunto de n vectores de \mathbb{R}^n , los cuales son ortogonales dos a dos, y cada uno tiene norma 1. Siempre es posible ortonormalizar vectores linealmente independientes para producir una base ortogonal, por ejemplo mediante el proceso de Gram Schmidt (Lay, 2012, Datta, 2010).

Los subespacios de un espacio vectorial V son subconjuntos propios de V que también son espacios vectoriales y por lo tanto disponen de bases para representar sus elementos.

Proyección en espacios de menor dimensión

Los grandes volúmenes de datos, generados por datos con alta dimensionalidad, son difíciles de analizar o de visualizar, por ejemplo, en el caso de las imágenes a colores, que son representadas por una matriz de píxeles, cada uno de los cuales representa 3 dimensiones, una para cada color (rojo, verde y azul), lo que significa que en alta resolución, la compresión y visualización son tareas que requieren de algoritmos eficientes; afortunadamente, con frecuencia, los datos pueden ser representados con tan solo algunas de las dimensiones. En particular, los procedimientos para comprimir imágenes utilizan técnicas para reducir su tamaño manteniendo la información más relevante o significativa. La reducción de dimensionalidad está fuertemente ligada al concepto de proyección ortogonal ya que se trata de proyectar la data en un subespacio de menor dimensión en el cual se puedan manejar los datos con mayor facilidad y extraer información relevante a menor costo.

Comencemos mostrando como se proyecta un vector x de \mathbb{R}^n en un subespacio de menor dimensión, por ejemplo una recta con vector director $v = (v_1, v_2, \dots, v_n) \in \mathbb{R}^n$. La recta es un subespacio de dimensión 1 que denotaremos con U . La proyección de x sobre U es un vector de U que denotaremos $Proy_U(x)$, la figura 1 ilustra el proceso.

De acuerdo con la definición, el vector proyección se puede escribir como combinación lineal de los vectores en la base del subespacio U , en este caso la base está formada solo por el vector v , por lo tanto $Proy_U(x) = \alpha v$, para algún escalar $\alpha \in \mathbb{R}$.

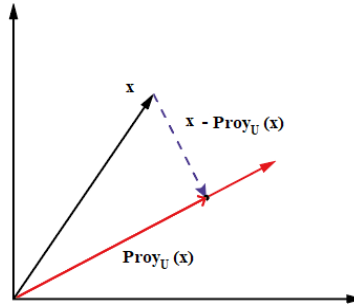


Fig. 1. Proyección en un subespacio de dimensión 1.

En este caso, la proyección es el vector de U más cercano a x , lo que significa que $\|x - \text{Proj}_U(x)\|$ es mínima y en consecuencia el segmento $x - \text{Proj}_U(x)$ es ortogonal a v , esto se caracteriza como:

$$\langle v, x - \text{Proj}_U(x) \rangle = \langle v, x - \alpha v \rangle = 0, \text{ para algún escalar } \alpha \in \mathbb{R}$$

Tomando en cuenta que el producto interno es bilineal, inferimos que

$$\langle v, x \rangle - \alpha \langle v, v \rangle = 0 \Leftrightarrow \alpha = \frac{\langle v, x \rangle}{\langle v, v \rangle}$$

Nótese que conociendo los vectores que generan el subespacio U , solo necesitamos el valor de α para obtener la proyección. En el ejemplo, si el vector v , que es la base del subespacio de proyección, tiene norma 1, entonces conseguir la proyección de x en U se reduce a hallar el escalar $\alpha = \langle v, x \rangle$, y posteriormente la proyección $\text{Proj}_U(x) = \alpha v$.

Si usamos el producto punto como producto interno tendremos que

$$\text{Proj}_U(x) = \left(\frac{v^T x}{v^T v} \right) v = \left(\frac{v^T x}{\|v\|^2} \right) v = v \frac{v^T x}{\|v\|^2} = \frac{v v^T}{\|v\|^2} v$$

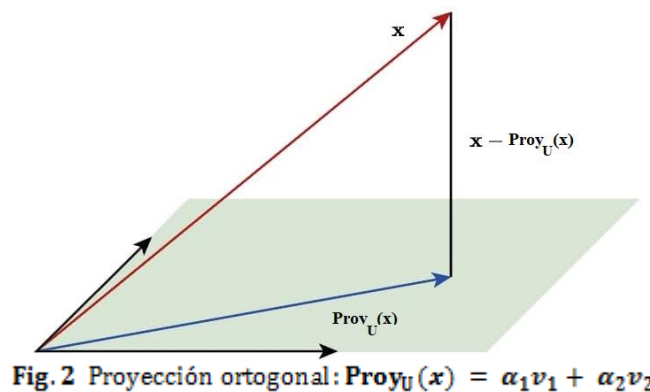
En esta última expresión notamos que el vector $\text{Proj}_U(x) = \left(\frac{v^T x}{v^T v} \right) v$ es calculado como el escalar $\alpha = \frac{v^T x}{\|v\|^2}$ multiplicado por el vector v ; sin embargo, esto es equivalente a calcularlo como el

producto de la matriz $P = \frac{v v^T}{\|v\|^2}$ por el vector v . A esta matriz P la llamamos matriz de proyección. Nótese que en este caso, sólo necesitamos al escalar α para representar la proyección, lo cual señala el camino hacia la reducción de la dimensionalidad.

Proyección en espacios de gran dimensión

Supongamos ahora que queremos proyectar un vector $\mathbf{x} = (x_1, x_2, x_3)$ de un espacio tridimensional en un subespacio U de dos dimensiones generado por los vectores \mathbf{v}_1 y \mathbf{v}_2 . Haciendo la analogía con la proyección en un subespacio de una dimensión, el vector proyección será una combinación lineal de \mathbf{v}_1 y \mathbf{v}_2 .

Adicionalmente, considerando que la proyección ortogonal nos dará la menor distancia, el vector $\mathbf{x} - \text{Proy}_U(\mathbf{x})$ será ortogonal al plano expandido por \mathbf{v}_1 y \mathbf{v}_2 , lo cual significa que es ortogonal a \mathbf{v}_1 y a \mathbf{v}_2 , ver figura 2.



En términos matriciales, consideremos B la matriz que tiene como columnas a \mathbf{v}_1 y \mathbf{v}_2 , y sea $\alpha = (\alpha_1, \alpha_2)$ el vector de escalares, entonces, las siguientes dos condiciones deben cumplirse:

- 1) $\text{Proy}_U(\mathbf{x}) = B\alpha$
- 2) $\langle \mathbf{x} - \text{Proy}_U(\mathbf{x}), \mathbf{v}_1 \rangle = 0$ y $\langle \mathbf{x} - \text{Proy}_U(\mathbf{x}), \mathbf{v}_2 \rangle = 0$

Sustituyendo (1) en (2) obtenemos

$$\begin{aligned} \langle \mathbf{x} - B\alpha, \mathbf{v}_1 \rangle &= 0 \quad \text{y} \quad \langle \mathbf{x} - B\alpha, \mathbf{v}_2 \rangle = 0 \\ \langle \mathbf{x}, \mathbf{v}_1 \rangle - \langle B\alpha, \mathbf{v}_1 \rangle &= 0 \quad \text{y} \quad \langle \mathbf{x}, \mathbf{v}_2 \rangle - \langle B\alpha, \mathbf{v}_2 \rangle = 0 \\ \mathbf{x}^T \mathbf{v}_1 - \alpha^T B^T \mathbf{v}_1 &= 0 \quad \text{y} \quad \mathbf{x}^T \mathbf{v}_2 - \alpha^T B^T \mathbf{v}_2 = 0 \\ \mathbf{x}^T [\mathbf{v}_1, \mathbf{v}_2] &= \alpha^T B^T [\mathbf{v}_1, \mathbf{v}_2] \equiv \alpha^T B^T B = \mathbf{x}^T B \end{aligned}$$

Una vez que los vectores v_1 y v_2 son linealmente independientes, la matriz $B^T B$ es cuadrada de orden 2 con columnas independientes, lo que significa que su inversa $(B^T B)^{-1}$ existe, y en consecuencia los escalares de la combinación pueden obtenerse como:

$$\alpha^T = x^T B (B^T B)^{-1} \equiv \alpha = (B^T B)^{-1} B^T x$$

Una vez que obtenemos α , la proyección será $\text{Proy}_U(x) = B\alpha = B(B^T B)^{-1} B^T x$.

A la matriz $P = B(B^T B)^{-1} B^T$ se le conoce como matriz de proyección. Si los vectores columnas de B son ortonormales, entonces $B^T B$ es la matriz identidad y el cálculo de la proyección se reduce a $\text{Proy}_U(x) = B B^T x$, siendo $P = B B^T$ la matriz de proyección. En este caso, para representar al vector $x = (x_1, x_2, x_3)$ de \mathbb{R}^3 en el subespacio U basta con tener $\alpha \in \mathbb{R}^2$, cumpliendo con el propósito de reducir la dimensionalidad.

En el caso general de la proyección de un vector n -dimensional x , en un subespacio U de dimensión $k \ll n$, del cual se conoce una base $B_k = \{v_1, v_2, v_3, \dots, v_k\}$ de vectores ortonormales, la matriz B tendrá como columnas los vectores $v_1, v_2, v_3, \dots, v_k$, y por lo tanto será de $n \times k$, la $\text{Proy}_U(x) = B B^T x$ será un vector n -dimensional que podrá ser representado por un vector $\alpha = B^T x$ que es k -dimensional.

Este resultado es la base que sustenta la reducción de dimensionalidad vía análisis de componentes principales.

Análisis de Componentes Principales y reducción de dimensionalidad

El análisis de componentes principales (PCA por sus siglas en inglés) es una de las técnicas más usadas para la comprensión y manejo de datos con altas dimensiones. Se ha utilizado por muchos años en el almacenamiento, compresión y visualización, de imágenes, aprovechando la característica de que generalmente las imágenes pueden ser representadas solo con algunas de las direcciones, y que muchas direcciones están altamente correlacionadas.

La siguiente gráfica ilustra la situación en dos dimensiones, en la cual los datos no están necesariamente en la misma recta; sin embargo su variación en un espacio de una dimensión es poco significativa, por lo que bien pudieran ser modelados por la recta que mejor los represente; este último término debe estar bien establecido, por ejemplo podríamos hablar de la recta de mejor ajuste en algún sentido; un modelo sugerido muchas veces es la recta que minimiza los errores cuadráticos en cada punto, sustentado por el método de los mínimos cuadrados lineales.

La figura 3 muestra el caso de variables linealmente correlacionadas, en el cual los datos bien pueden ser modelados por una recta.

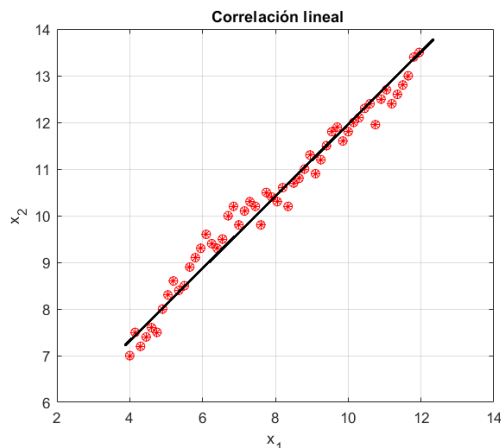


Fig. 3. Correlación lineal entre variables

Supongamos que tenemos un conjunto de datos \mathbf{X} que contiene N vectores $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N$, cada uno de ellos representando d características, es decir cada uno es d -dimensional. Hablamos entonces de una muestra de N objetos en un espacio d -dimensional que quisiéramos representar en un espacio k -dimensional de menor dimensión ($k \ll d$) de tal forma que esta representación sea lo más similar posible a la muestra original \mathbf{X} . El objetivo del análisis de componentes principales es minimizar el error de reconstrucción promedio de las proyecciones ortogonales con respecto a la data original.

Ahora bien, sea $\mathbf{B}_d = \{\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \dots, \mathbf{v}_d\}$ una base ortonormal de un espacio d -dimensional, entonces cada elemento de \mathbf{X}_i , del conjunto de datos, se puede escribir como una combinación lineal de los vectores de la base \mathbf{B}_d , esto es, existe un vector de escalares $\alpha = (\alpha_{i1}, \alpha_{i2}, \dots, \alpha_{id})$ tal que:

$$\mathbf{X}_i = \alpha_{i1}\mathbf{v}_1 + \alpha_{i2}\mathbf{v}_2 + \dots + \alpha_{id}\mathbf{v}_d$$

Una vez que $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d$ son ortonormales $\mathbf{v}_j^T \mathbf{X}_i = \alpha_{ij}$, con $j = 1, \dots, d$, lo cual significa que α_{ij} representa la proyección de \mathbf{X}_i sobre el subespacio unidimensional generado por \mathbf{v}_j . Si ahora tomamos dos vectores \mathbf{v}_j y \mathbf{v}_p de \mathbf{B}_d , obtendremos que $[\mathbf{v}_j, \mathbf{v}_p]^T \mathbf{X}_i = [\alpha_{ij}, \alpha_{ip}]$ es una representación de \mathbf{X}_i en un espacio de dos dimensiones. Es claro que con este procedimiento

podríamos obtener siempre una representación de los datos en un espacio de menor dimensión, y también pudiéramos reconocer que para un cierto $k \ll n$ la representación es adecuada.

Con todas estas consideraciones, podríamos seleccionar k vectores de \mathbf{B}_d ($k \ll d$) y definir \mathbf{B} como la matriz cuyas columnas son los k vectores ortonormales seleccionados, y la proyección del dato \mathbf{X}_i sobre el subespacio U generado por estos k vectores se calcula como $\text{Proy}_U(\mathbf{X}_i) = \mathbf{B}\mathbf{B}^T\mathbf{X}_i$, y puede ser representada por un vector de coordenadas (reconocido como code) $\alpha_i = (\alpha_{i1}, \alpha_{i2}, \dots, \alpha_{ik}) = \mathbf{B}^T\mathbf{X}_i$.

Como hemos dicho anteriormente, la proyección $\text{Proy}_U(\mathbf{X}_i)$ es un elemento del espacio d -dimensional, y por lo tanto se puede expresar como combinación lineal de los vectores de \mathbf{B}_d , $\text{Proy}_U(\mathbf{X}_i) = \beta_1\mathbf{v}_1 + \beta_2\mathbf{v}_2 + \dots + \beta_d\mathbf{v}_d$.

Sin pérdida de generalidad, asúmanos que los k vectores que generan el subespacio U son los vectores $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$ de \mathbf{B}_d , los cuales son las columnas de la matriz \mathbf{B} ; entonces

$$\begin{aligned} \text{Proy}_U(\mathbf{X}_i) &= \beta_1\mathbf{v}_1 + \beta_2\mathbf{v}_2 + \beta_k\mathbf{v}_k + \beta_{k+1}\mathbf{v}_{k+1} \dots + \beta_d\mathbf{v}_d \\ &= \sum_{n=1}^k \beta_n\mathbf{v}_n + \sum_{n=k+1}^d \beta_n\mathbf{v}_n = \tilde{\mathbf{X}}_i + \sum_{n=k+1}^d \beta_n\mathbf{v}_n. \end{aligned}$$

La primera suma es $\tilde{\mathbf{X}}_i = \mathbf{B}\alpha$, la aproximación a los datos en el subespacio U , con $\alpha = (\beta_1, \beta_2, \dots, \beta_k)$, y la segunda suma es un vector en el complemento ortogonal de U . Lo que se plantea en el análisis de componentes principales es descartar la segunda suma, tal como lo hemos indicado previamente, y definir al subespacio generado por los vectores $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$ como el subespacio principal, aceptando al vector $\tilde{\mathbf{X}}_i$ como la representación o aproximación de \mathbf{X}_i en el subespacio U de menor dimensión.

Considerando N objetos en el conjunto de datos $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N$ el objetivo entonces es hallar el vector de coordenadas α y los vectores \mathbf{v}_i , con $i = 1, \dots, k$, tal que el error cuadrático de reconstrucción promedio se minimice. Si denotamos \mathbf{E}_k este error de reconstrucción cuando deseamos representar la data a un espacio k -dimensional, el problema consiste en:

$$\text{minimizar } \mathbf{E}_k = \frac{1}{N} \sum_{n=1}^N \|\mathbf{X}_n - \tilde{\mathbf{X}}_n\|^2 = \frac{1}{N} \sum_{n=1}^N (\mathbf{X}_n - \tilde{\mathbf{X}}_n)^T (\mathbf{X}_n - \tilde{\mathbf{X}}_n)$$

Con el propósito de facilitar la notación asumimos que la esperanza del conjunto de datos es cero y que disponemos de una base \mathbf{B}_d de vectores ortonormales en el espacio d -dimensional donde se encuentran los datos, y que para todo n queremos hallar $\tilde{\mathbf{X}}_n = \mathbf{B}\alpha_n$ similar a \mathbf{X}_n .

La medida de similaridad que se plantea es el cuadrado de la distancia euclídea $\|\mathbf{X}_n - \tilde{\mathbf{X}}_n\|^2$, siendo el objetivo minimizar el promedio de los errores cuadráticos o el error de reconstrucción (Pearson, 1901).

Se plantea entonces un problema de optimización, y para hallar el $\tilde{\mathbf{X}}_n$ óptimo, para un n dado $n = 1, \dots, N$, debemos encontrar el conjunto de escalares α_n y el conjunto de vectores $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$ óptimos, por lo tanto se trata de un problema de optimización multivariable. Una forma simple de lograr este objetivo es hallar primero el conjunto de escalares óptimos dada una base fija de vectores $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$, y luego hallar los vectores óptimos, que conformarán el subespacio principal.

Una vez que hemos escogido la norma euclídea, como la herramienta para medir la similaridad, el conjunto de escalares óptimos, para un $\mathbf{X}_n = (X_{n1}, \dots, X_{nd})$ dado, es aquel que anula las derivadas parciales de E_k con respecto a cada coordenada en α_n , las cuales son funciones de $\tilde{\mathbf{X}}_n$, por lo tanto debemos aplicar la regla de la cadena.

$$\frac{\partial E_k}{\partial \alpha_{in}} = \frac{\partial E_k}{\partial \tilde{\mathbf{X}}_n} \frac{\partial \tilde{\mathbf{X}}_n}{\partial \alpha_{in}}$$

Considerando que el error de reconstrucción viene dado por $E_k = \frac{1}{N} \sum_{n=1}^N \|\mathbf{X}_n - \tilde{\mathbf{X}}_n\|^2$, tenemos que para un dato particular \mathbf{X}_n el término correspondiente en la sumatoria es:

$$\begin{aligned} \|\mathbf{X}_n - \tilde{\mathbf{X}}_n\|^2 &= \left\| \begin{pmatrix} X_{n1} \\ X_{n2} \\ \vdots \\ X_{nd} \end{pmatrix} - \begin{pmatrix} \tilde{X}_{n1} \\ \tilde{X}_{n2} \\ \vdots \\ \tilde{X}_{nd} \end{pmatrix} \right\|^2 = \left\| \begin{pmatrix} X_{n1} - \tilde{X}_{n1} \\ X_{n2} - \tilde{X}_{n2} \\ \vdots \\ X_{nd} - \tilde{X}_{nd} \end{pmatrix} \right\|^2 \\ &= (X_{n1} - \tilde{X}_{n1})^2 + (X_{n2} - \tilde{X}_{n2})^2 + \dots + (X_{nd} - \tilde{X}_{nd})^2 \end{aligned}$$

Por lo tanto,

$$\begin{aligned} \frac{\partial E_k}{\partial \tilde{\mathbf{X}}_n} &= \frac{1}{N} \left(-2(X_{n1} - \tilde{X}_{n1}) - 2(X_{n2} - \tilde{X}_{n2}) - \dots - 2(X_{nd} - \tilde{X}_{nd}) \right) \\ &= \frac{-2}{N} (X_n - \tilde{\mathbf{X}}_n)^T \\ \Rightarrow \frac{\partial E_k}{\partial \tilde{\mathbf{X}}_n} &= \frac{1}{N} \left(-2(X_{n1} - \tilde{X}_{n1}) - 2(X_{n2} - \tilde{X}_{n2}) - \dots - 2(X_{nd} - \tilde{X}_{nd}) \right) \\ &= \frac{-2}{N} (X_n - \tilde{\mathbf{X}}_n)^T \end{aligned}$$

Ahora bien, una vez que $\tilde{\mathbf{X}}_n = \alpha_{1n} \mathbf{v}_1 + \alpha_{2n} \mathbf{v}_2 + \dots + \alpha_{kn} \mathbf{v}_k$ se deduce que:

$$\frac{\partial \tilde{X}_n}{\partial \alpha_{in}} = v_i$$

Por lo tanto, el valor óptimo para el parámetro α_{in} viene dado por:

$$\begin{aligned} \frac{\partial E_k}{\partial \alpha_{in}} &= \frac{\partial E_k}{\partial \tilde{X}_n} \frac{\partial \tilde{X}_n}{\partial \alpha_{in}} = -\frac{2}{N} (X_n - \tilde{X}_n)^T v_i = -\frac{2}{N} (X_n^T v_i - \tilde{X}_n^T v_i) = 0 \\ &\Rightarrow X_n^T v_i - \alpha_{in} = 0 \Rightarrow \alpha_{in} = X_n^T v_i \end{aligned} \quad (1)$$

De esta manera, para cada X_n ($n = 1, \dots, N$) podemos obtener α_{in} ($i = 1, \dots, k$), reconociendo que el valor óptimo de α_{in} coincide con aquel que nos da la proyección ortogonal del elemento X_n de la data original en el subespacio 1-dimensional generado por el vector v_i . La generalización de este resultado apoya la propuesta inicial de considerar la proyección ortogonal $\text{Proj}_U(X) = \tilde{X}$ como una representación de la de la data original d -dimensional X en el subespacio k -dimensional U ($k \ll d$) producido por k -vectores dados. El próximo paso será la escogencia de los vectores v_1, v_2, \dots, v_k que definen el subespacio principal.

Tal como quedó establecido en el resultado (1) cada coordenada del vector de escalares para un determinado X_n se calcula como $\alpha_{in} = X_n^T v_i$, y la proyección se expresa en términos de la combinación lineal $\tilde{X}_n = \alpha_{1n} v_1 + \alpha_{2n} v_2 + \dots + \alpha_{kn} v_k$, por lo tanto,

$$\begin{aligned} \tilde{X}_n &= \sum_{i=1}^k \alpha_{in} v_i = \sum_{i=1}^k (X_n^T v_i) v_i = \left(\sum_{i=1}^k v_i v_i^T \right) X_n = BB^T X_n \\ \tilde{X}_n &= \sum_{i=1}^k \alpha_{in} v_i = \sum_{i=1}^k (X_n^T v_i) v_i = \sum_{i=1}^k (v_i^T X_n) v_i \end{aligned} \quad (2)$$

La última igualdad en (2) se justifica por la simetría del producto interno.

De igual manera, X_n se expresa como combinación lineal de todos los vectores de la base B_d , lo cual podemos expresar en términos de dos sumatorias, esto es:

$$\begin{aligned} X_n &= \sum_{i=1}^d \alpha_{in} v_i = \sum_{i=1}^d (v_i^T X_n) v_i = \sum_{i=1}^k (v_i^T X_n) v_i + \sum_{i=k+1}^d (v_i^T X_n) v_i \\ &\Rightarrow X_n = \sum_{i=1}^d \alpha_{in} v_i = \tilde{X}_n + \sum_{i=k+1}^d (v_i^T X_n) v_i \\ &\Rightarrow X_n - \tilde{X}_n = \sum_{i=k+1}^d (v_i^T X_n) v_i \Rightarrow \|X_n - \tilde{X}_n\|^2 = \left\| \sum_{i=k+1}^d (v_i^T X_n) v_i \right\|^2 \end{aligned} \quad (3)$$

Por lo tanto, el error al aproximar X_n por \tilde{X}_n , en la medida establecida, es igual a la aproximación que resultaría si usáramos el complemento ortogonal del subespacio principal, y este hecho puede ser usado equivalentemente para hallar la solución óptima, tal como veremos a continuación.

Veamos como reformular la expresión $E_k = \frac{1}{N} \sum_{n=1}^N \|X_n - \tilde{X}_n\|^2$ del error cuadrático promedio que queremos optimizar usando el resultado en (3).

$$E_k = \frac{1}{N} \sum_{n=1}^N \|X_n - \tilde{X}_n\|^2 = \frac{1}{N} \sum_{n=1}^N \left\| \sum_{i=k+1}^d (v_i^T X_n) v_i \right\|^2$$

En esta última expresión podemos identificar escalares $\gamma_i = (v_i^T X_n)$, $i = k+1, \dots, d$ que intervienen en una combinación lineal de los vectores $v_{k+1}, v_{k+2}, \dots, v_d$ del complemento ortogonal.

$$\sum_{i=k+1}^d (v_i^T X_n) v_i = \gamma_{k+1} v_{k+1} + \gamma_{k+2} v_{k+2} + \dots + \gamma_d v_d$$

Por lo tanto, la norma euclídea al cuadrado no es más que el producto escalar de un vector del complemento ortogonal por si mismo.

$$(\gamma_{k+1} v_{k+1} + \gamma_{k+2} v_{k+2} + \dots + \gamma_d v_d)^T (\gamma_{k+1} v_{k+1} + \gamma_{k+2} v_{k+2} + \dots + \gamma_d v_d)$$

En donde $v_i^T v_j = 0$ si $i \neq j$ y $v_i^T v_j = 1$ si $i = j$, $\forall i, j = k+1, \dots, d$, por ser estos vectores ortonormales. De esta manera

$$\left\| \sum_{i=k+1}^d (v_i^T X_n) v_i \right\|^2 = \sum_{i=k+1}^d (v_i^T X_n)^2$$

Luego, sustituyendo este resultado en el objetivo a optimizar, tenemos que:

$$\begin{aligned} E_k &= \frac{1}{N} \sum_{n=1}^N \left\| \sum_{i=k+1}^d (v_i^T X_n) v_i \right\|^2 = \frac{1}{N} \sum_{n=1}^N \sum_{i=k+1}^d (v_i^T X_n)^2 \\ \Rightarrow E_k &= \frac{1}{N} \sum_{n=1}^N \sum_{i=k+1}^d v_i^T X_n v_i^T X_n = \frac{1}{N} \sum_{n=1}^N \sum_{i=k+1}^d v_i^T X_n X_n^T v_i^T \end{aligned}$$

El último resultado, que obedece a la simetría de del producto interno, y una vez que las sumas son intercambiables, tenemos que:

$$E_k = \sum_{i=k+1}^d v_i^T \left(\frac{1}{N} \sum_{n=1}^N X_n X_n^T \right) v_i = \sum_{i=k+1}^d v_i^T S v_i$$

Podemos identificar en esta nueva expresión del error cuadrático de reconstrucción a la matriz S como la matriz de covarianzas de datos centrados.

Con esta reformulación de E_k , el problema se traduce en hallar los vectores v_j , con j en $\{k+1, k+2, \dots, d\}$ del complemento ortogonal al subespacio principal, sujeto a que estos vectores sean ortonormales. Esto lo podemos escribir como:

$$\text{minimizar } E_k = v_{k+1}^T S v_{k+1} + v_{k+2}^T S v_{k+2} + \dots + v_d^T S v_d$$

$$\text{sujeto a: } v_i^T v_j = \delta_{ij} \quad \text{con } i, j = \{k+1, k+2, \dots, d\}.$$

Podemos resolver este problema usando el método de los multiplicadores de Lagrange (ver Thomas, 2005). La función Lagrangiana sería:

$$L = v_{k+1}^T S v_{k+1} + \dots + v_d^T S v_d + \lambda_{k+1} (1 - v_{k+1}^T v_{k+1}) + \dots + \lambda_d (1 - v_d^T v_d).$$

La solución se hallaría resolviendo el sistema de Lagrange, para lo cual se igualan a cero las derivadas parciales con respecto a cada multiplicador λ_j , $j = k+1, \dots, d$, y también las derivadas parciales con respecto a los vectores v_j , $j = k+1, \dots, d$.

$$\frac{\partial L}{\partial \lambda_j} = 1 - v_j^T v_j = 0 \Leftrightarrow v_j^T v_j = 1$$

$$\frac{\partial L}{\partial v_j} = 2v_j^T S - 2\lambda v_j^T = 0 \Leftrightarrow S v_j = \lambda v_j$$

Finalmente, podemos notar, que los vectores que optimizan el error cuadrático promedio son autovectores de la matriz de covarianzas. En este caso, el valor del objetivo estaría dado por:

$$\begin{aligned} E_k &= v_{k+1}^T \lambda_{k+1} v_{k+1} + v_{k+2}^T \lambda_{k+2} v_{k+2} + \dots + v_d^T \lambda_d v_d \\ \Rightarrow E_k &= \lambda_{k+1} v_{k+1}^T v_{k+1} + \lambda_{k+2} v_{k+2}^T v_{k+2} + \dots + \lambda_d v_d^T v_d \\ \Rightarrow E_k &= \lambda_{k+1} + \lambda_{k+2} + \dots + \lambda_d \end{aligned}$$

Así, para saber el valor mínimo de E_k basta con hallar los $d-k$ autovalores de menor valor de S , que por ser una matriz simétrica y semidefinida positiva tiene autovalores reales no negativos. Sin embargo, los correspondientes autovectores $v_{k+1}, v_{k+2}, \dots, v_d$ pertenecen al complemento ortogonal del subespacio principal. Por lo tanto, y los vectores que necesitamos en la base B son los autovectores correspondientes a los k autovalores de mayor valor.

Esto puede apreciarse en la figura 4, que presenta una data en dos dimensiones y señala los dos autovectores de la matriz de covarianzas; se puede observar que el autovector asociado al mayor

autovalor (resaltado en rojo) representará la mejor base para la proyección, mientras que el autovector asociado al autovalor de módulo mínimo sugiere la magnitud del error.

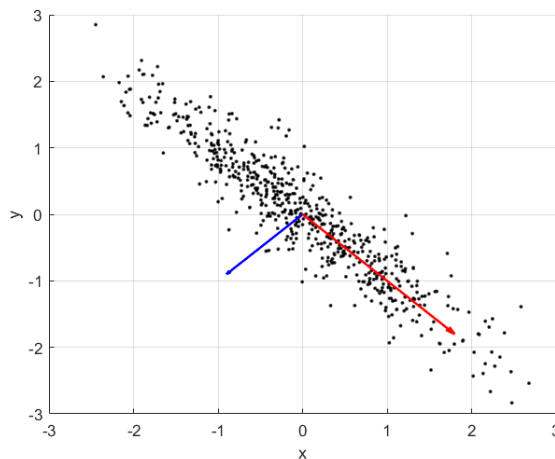


Fig. 4. Representación del subespacio principal

Algoritmo básico (reducción de dimensionalidad vía PCA)

Una vez que entendemos cómo obtener la base del espacio de proyección, podemos retomar el cálculo de los escalares (el code) de la combinación lineal que permite construir la proyección para cada dato \mathbf{X}_n , con $n = 1, \dots, N$. Esto nos habilita para definir un algoritmo con los pasos básicos del proceso.

En la solución del problema hemos asumido que la data es centrada, o que la media es cero; solo para disminuir la dificultad en la derivación, pero obtendríamos el mismo resultado con cualquier valor de la media; aunque por cuestiones numérico-computacionales, es preferible restar la media a los datos en este tipo de análisis. Para mayor información sobre estas prácticas y otras técnicas eficientes para implementaciones de PCA se sugiere ver Deisenroth 2020, y James, 2013)

Otra práctica conveniente es normalizar la data; y esto lo hacemos dividiendo en cada dirección por la desviación estándar, lo cual libera la dependencia de unidades y garantiza una varianza unitaria en cada dimensión sin alterar la correlación entre variables.

En resumen, los primeros pasos del algoritmo consisten en centrar y normalizar la data; eso es, por cada punto \mathbf{X}_n :

$$X_n^{(*)} \leftarrow \frac{X_n - \mu}{\sigma}$$

Luego, debemos hallar la matriz de covarianzas S de la data $\mathbf{X}^* = \{\mathbf{X}_1^*, \dots, \mathbf{X}_n^*\}$, y calcular los k autovalores de mayor valor. Los k autovectores asociados a estos autovalores formarán las columnas de la matriz de proyección, con la cual finalmente obtenemos la data proyectada.

$$\tilde{\mathbf{X}}_n = \text{Proy}_U(\mathbf{X}_n^*) = \mathbf{B}\mathbf{B}^T(\mathbf{X}_n^*)$$

Implementación en Matlab

A continuación, se presenta un código básico para la reducción de dimensionalidad; en el cual se crea una data aleatoria, que será sujeta a la reducción de dos dimensiones a una dimensión. La salida es el vector de escalares alfa (en este caso un escalar, ya que el subespacio de proyección es unidimensional).

```

1 function [alfa]=PCAO(X)
2 % Dimensionality reduction via PCA
3 % This code takes a random dataset X, with N two-dimensional
4 % vectors, and project it onto a one-dimensional subspace, to
5 % illustrate the Principal Component Analysis steps.
6
7 N=length(X);
8
9 % Centering and normalizing
10 mu1=mean(X(:,1)); s1=std2(X(:,1));
11 mu2=mean(X(:,2)); s2=std2(X(:,2));
12 X(:,1)=(X(:,1)-mu1*ones(N,1))/s1;
13 X(:,2)=(X(:,2)-mu2*ones(N,1))/s2;
14
15 % Finding the covariance matrix
16 CovX=cov(X);
17
18 % Identifying the principal component
19 [V,E]=eig(CovX);
20 e = diag(E);
21 [me,ind]=max(e);
22
23 % Defining the projection subspace
24 u=V(:,ind);
25
26 % Finding the optimal scalars(code)
27 alfa= X*u;
28
29 % Calculating the projection (Decode)
30 Proj_u=alfa*u';
31
32 % Visualizing results
33 figure;
34 scatter(X(:,1), X(:,2), 'k. ');
35 hold on;
36 scatter(Proj_u(:,1), Proj_u(:,2), 'g. ');
37 title('PCA: Reduccion de dimensionalidad');
38 xlabel('x');
39 ylabel('y');
40 end

```

Fig. 5. Código básico en Matlab para PCA

La figura 6 muestra los resultados de la ejecución del código (Fig. 5) para un conjunto $N=500$ datos aleatorios.

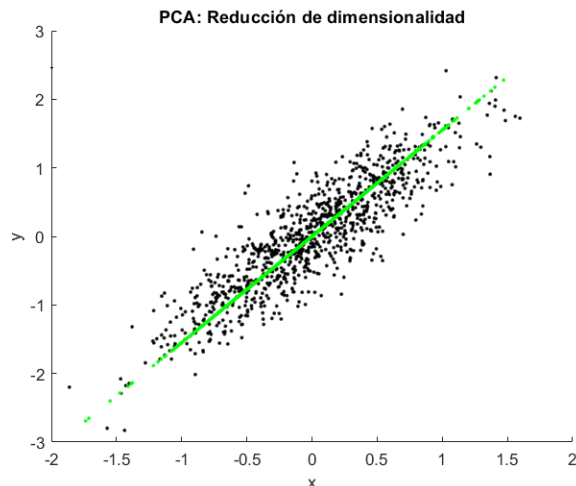


Fig. 6. PCA para reducción de datos de 2 dimensiones a una dimensión

Esta filosofía puede extenderse al caso general de N datos con d características que se desee proyectar en un espacio de $k \ll d$ dimensiones; así, seleccionaríamos k autovectores asociados a los k autovalores de mayor valor de la matriz de covarianzas.

Nótese que en el paso 20 del código estamos calculando la descomposición espectral de la matriz que representa a los datos; sin embargo, lo conveniente en la práctica es calcular solo k autovectores asociados a los autovalores de mayor valor, y esto puede hacerse con métodos de proyección para el cálculo de autovalores; lo que sería recomendable en el caso general de aplicaciones con alta dimensionalidad. Para una explicación detallada de estas técnicas se recomienda ver (Trefethen, 1997). Una vez que la matriz S es simétrica y deficiente en rango, una propuesta con cálculo estable sería hallar la descomposición en valores singulares truncada de la matriz X . El tópico ha sido tratado ampliamente en bibliografías recientes, ver por ejemplo (Deisenroth, 2020).

Conclusiones

Bajo la suposición de tener un conjunto de N datos X_1, X_2, \dots, X_N con media $\mathbf{0}$, en un espacio d -dimensional, y disponer de una base ortogonal de vectores para este espacio, se han descrito en

forma detallada los elementos de la matemática que sirven de base para uno de los algoritmos más usado en la reducción de dimensionalidad, como lo es el análisis de componentes principales o PCA. En forma incremental se describieron las bases del algoritmo que soluciona el problema. El análisis de componentes principales ha sido por muchos años la referencia para hacer reducción de dimensionalidad y existen muchas referencias que lo definen y documentan, generalmente desde perspectivas de alto nivel de especialización en estadística, o bajo el enfoque de nuevos paradigmas como el aprendizaje automático (o Machine Learning) sin embargo, las herramientas básicas para entender el funcionamiento del algoritmo y programarlo en el lenguaje de programación de preferencia, yacen en el álgebra lineal, por lo que puede ser considerado como una aplicación del álgebra lineal.

Otro esquema de derivación del algoritmo consiste en maximizar la varianza en el subespacio de proyección para retener información de interés (Deisenroth, 2020). El enfoque es elegante pero requiere de elementos más sofisticados ligados a la estadística. En este trabajo hemos querido usar herramientas más básicas de la matemática.

Describimos normas relacionadas con diferentes productos internos, aunque usamos el producto escalar clásico en la derivación del algoritmo, por lo cual queda abierta la posibilidad de generar y comparar con otros algoritmos basados en diferentes normas. De igual manera trabajos para el cálculo efectivo y de bajo costo computacional de los autovalores de módulo máximo de la matriz de covarianza se proponen a futuro.

Referencias

1. Datta, B.N (2010). *Numerical Linear Algebra and Applications*. Society for Industrial and Applied Mathematics; 2a. Edición.
2. De la Puente, V., C. (2018). *Estadística descriptiva e inferencial*, Ediciones IDT CB, Madrid, España.
3. Deisenroth, M.P., Faisal, A., Soon, C. (2020). *Mathematics for Machine Learning*. Cambridge University Press. <https://mml-book.com>.
4. James, G., Witten, D., Hastie, T., Tibshirani, R. (2013). *An introduction to Statistical learning with applications in R*, Springer, New York.
5. Lay, D. (2012). *Álgebra lineal y sus aplicaciones*. Pearson Education, México.

6. Pearson, K. (1901). *On lines and planes of closest fit to systems of points in space*. Philosophical Magazine, 2(6), 559 – 572
7. Rencher, A.C., Bruce Schaalje G.B. (1996). (2008). *Linear models in statistics*. Hoboken, N.J: Wiley-Interscience, 6a. Edición.
8. Thomas, G.B. (2005). *Cálculo. Varias Variables*, Pearson, Addison Wesley, 11a. Edición.
9. Salazar, P., C., Del Castillo, G., S. (2017). *Fundamentos Básicos de Estadística*, <http://www.dspace.uce.edu.ec/handle/25000/13720>, 1ª. Edición
10. Trefethen, L. N. & Bau III, D. (1997). *Numerical linear algebra*. Siam, Philadelphia.