

Técnicas de Reducción de Dimensión

A. M. Alvarez-Meza, Ph.D.

D. F. Collazos-Huertas, Ph.D(c)

amalvarezme@unal.edu.co, dfcollazos@unal.edu.co

Minería de datos

Departamento de Matemáticas y Estadística

Universidad Nacional de Colombia-sede Manizales



Contenido

- 1 Definición de reducción de dimensión
- 2 La maldición de la dimensionalidad
- 3 Enfoques principales para la DR
- 4 Análisis de componentes principales (PCA)
- 5 Locally Linear Embedding
- 6 Otras técnicas de reducción de dimensión

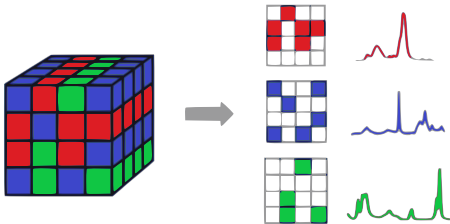
- 1 Definición de reducción de dimensión
- 2 La maldición de la dimensionalidad
- 3 Enfoques principales para la DR
- 4 Análisis de componentes principales (PCA)
- 5 Locally Linear Embedding
- 6 Otras técnicas de reducción de dimensión

Reducción de dimensión

Muchos problemas de **Machine Learning** implican **millones de características** para cada instancia de entrenamiento.

- Extremadamente lento
- Es mucho más difícil encontrar una buena solución.

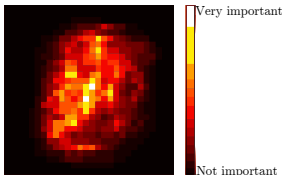
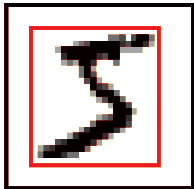
Este problema a menudo se conoce como *curse of dimensionality*.



En el mundo real, es posible reducir considerablemente el número de características, convirtiendo un “**intractable problem**” en un “**tractable problem**”.

Reducción de dimensión: Ejemplo

Consideremos las imágenes de la base de datos **MNIST**:

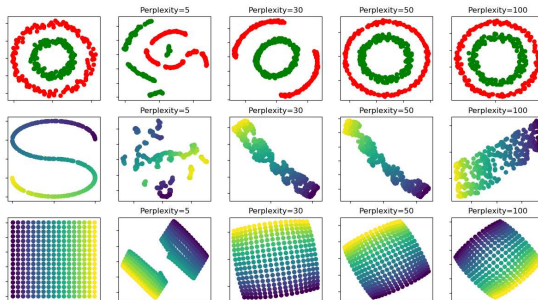


- Los píxeles en los bordes de la imagen **casi siempre son blancos**, por lo que puede eliminar estos píxeles **sin perder mucha información**.
- La imagen de relevancia, confirma que **estos píxeles no son importantes para la tarea de clasificación**.

En el proceso de reducción de la dimensionalidad **se pierde algo de información**.

Reducción de dimensión: TIPS

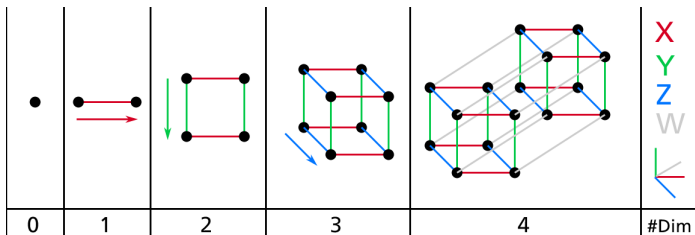
- En algunos casos, reducir la dimensionalidad de los datos de entrenamiento **puede filtrar algunos ruidos y detalles innecesarios** y, por lo tanto, **generar un mayor rendimiento**.
- Además de acelerar el entrenamiento, la reducción de dimensionalidad también **es extremadamente útil para la visualización de datos**.



- 1 Definición de reducción de dimensión
- 2 La maldición de la dimensionalidad**
- 3 Enfoques principales para la DR
- 4 Análisis de componentes principales (PCA)
- 5 Locally Linear Embedding
- 6 Otras técnicas de reducción de dimensión

Curse of Dimensionality

Estamos tan acostumbrados a vivir en **tres dimensiones** que nuestra intuición nos **falla** cuando tratamos de imaginar **un espacio de alta dimensión**.



Resulta que muchas cosas se comportan de manera muy diferente en el **espacio de alta dimensión**

Curse of Dimensionality: Ejemplo

- En un hipercubo de unidad de 10.000 dimensiones, **la mayoría de los puntos en un hipercubo de alta dimensión están muy cerca del borde.**
- En un **cuadrado de la unidad**, la distancia entre estos dos puntos será, en promedio, aproximadamente 0.52. Si elige dos puntos aleatorios **en un cubo 3D de unidad**, la distancia promedio será aproximadamente 0.66.
- ¿Qué pasa con dos puntos elegidos al azar en un hipercubo de un millón de dimensiones? ¡Bueno, la distancia promedio, será de 408.25. Es decir, que la mayoría de las instancias de entrenamiento estén muy lejos unas de otras.

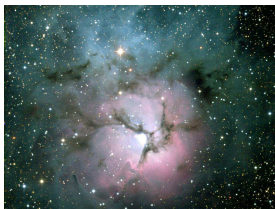
Cuanto más dimensiones tenga el conjunto de entrenamiento, mayor será el riesgo de **“over-fitting”**.

Curse of Dimensionality: Solución

Solución a la maldición de la dimensionalidad: aumentar el tamaño del conjunto de entrenamiento para alcanzar una densidad suficiente de instancias de entrenamiento.

Sin embargo, el número de instancias de entrenamiento requeridas para alcanzar una densidad dada crece exponencialmente con el número de dimensiones.

- Con solo 100 características, necesitaría más instancias que átomos en el universo.

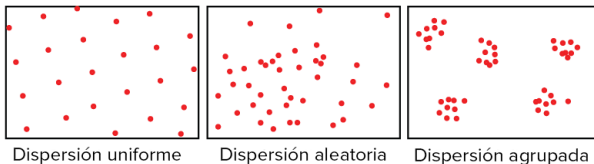


Contenido

- 1 Definición de reducción de dimensión
- 2 La maldición de la dimensionalidad
- 3 Enfoques principales para la DR**
- 4 Análisis de componentes principales (PCA)
- 5 Locally Linear Embedding
- 6 Otras técnicas de reducción de dimensión

Enfoques principales: Proyección

- En los problemas del mundo real, las instancias de entrenamiento no se distribuyen de manera uniforme en todas las dimensiones.

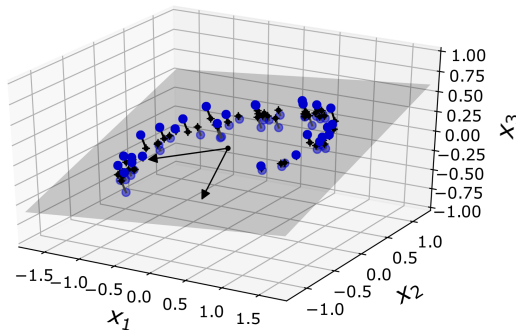


- Muchas características **son casi constantes**, mientras que otras están **altamente correlacionadas**.

Todas las instancias de entrenamiento realmente se encuentran dentro (o cerca de) un subespacio de dimensiones mucho más bajas.

Proyección: Ejemplo I

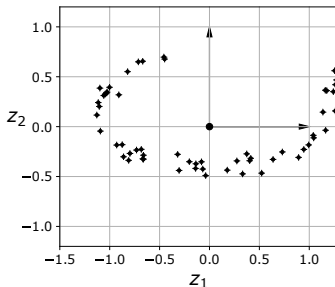
En la Figura se puede ver un conjunto de datos en 3D representado por círculos.



Todas las instancias de entrenamiento se encuentran cerca de un plano: este es un subespacio de menor dimensión (2D) del espacio de alta dimensión (3D).

Proyección: Ejemplo I

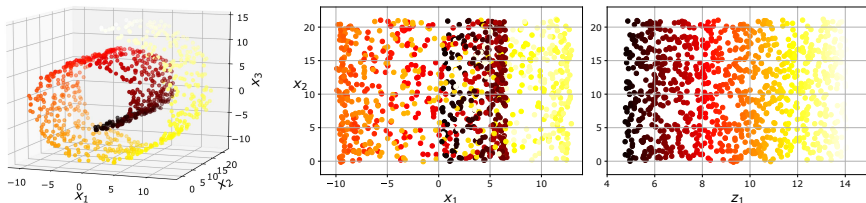
Si **proyectamos cada instancia de entrenamiento perpendicularmente en este subespacio**, obtenemos el nuevo conjunto de datos 2D:



- Tengamos en cuenta que los ejes corresponden a las nuevas características z_1 y z_2 .
- Sin embargo, la proyección no siempre es el mejor enfoque para la reducción de la dimensionalidad.

Proyección: Ejemplo II

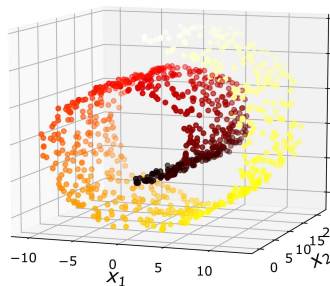
En muchos casos, el subespacio puede girar y girar, como en el famoso conjunto de datos “*Swiss roll*”:



- Proyectar en un plano (por ejemplo, quitando x_3) **aplastaría diferentes capas del rollo suizo** (Centro).
- Lo que realmente deseamos es **desenrollar el rollo suizo para obtener el conjunto de datos 2D** (Derecha).

Enfoques principales: Manifold learning

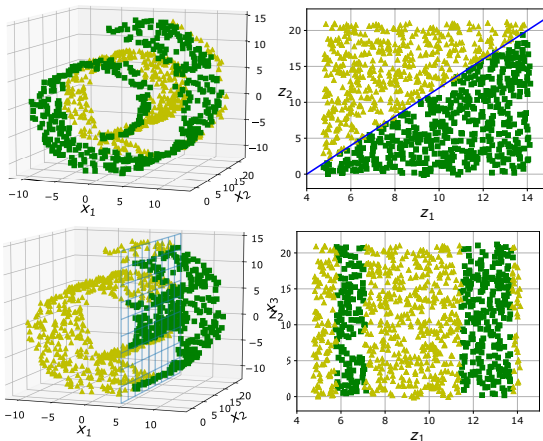
- El “*Swiss roll*” es un ejemplo de una variedad 2D.
- Un múltiple 2D es una forma 2D que se puede doblar y torcer en un espacio de dimensiones superiores.
- Una variedad d -dimensional es una parte de un espacio n -dimensional (**donde** $d < n$) que localmente se parece a un hiperplano d -dimensional.



Manifold assumption: la mayoría de los conjuntos de datos de alta dimensión del mundo real se encuentran cerca de un manifold de mucha más baja dimensión.

Manifold learning: Ejemplo I

El **Implicit assumption**: la tarea en cuestión (por ejemplo, clasificación o regresión) será más simple si se expresa en el espacio de dimensiones inferiores del manifold.



Manifold learning: Ejemplo I

- En la **fila superior de la Figura**, el rollo suizo se divide en dos clases: en el espacio 3D (a la izquierda), **el límite de decisión sería bastante complejo**, pero en el espacio múltiple desenrollado 2D (en el derecha), **el límite de decisión es una simple línea recta**.
- De lo contrario, en la **fila inferior de la Figura**, el límite de decisión se encuentra en $x_1 = 5$. **Este límite de decisión parece muy simple en el espacio 3D original**, pero **se ve más complejo en el múltiple desenrollado**.

Si se reduce la dimensionalidad antes de entrenar un modelo, acelerará el entrenamiento, pero **no siempre** puede conducir a una solución mejor.

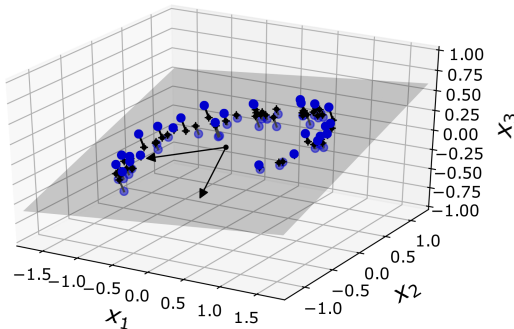
Contenido

- 1 Definición de reducción de dimensión
- 2 La maldición de la dimensionalidad
- 3 Enfoques principales para la DR
- 4 Análisis de componentes principales (PCA)**
- 5 Locally Linear Embedding
- 6 Otras técnicas de reducción de dimensión

PCA es el algoritmo de reducción de dimensionalidad más popular.

¿Cómo funciona?:

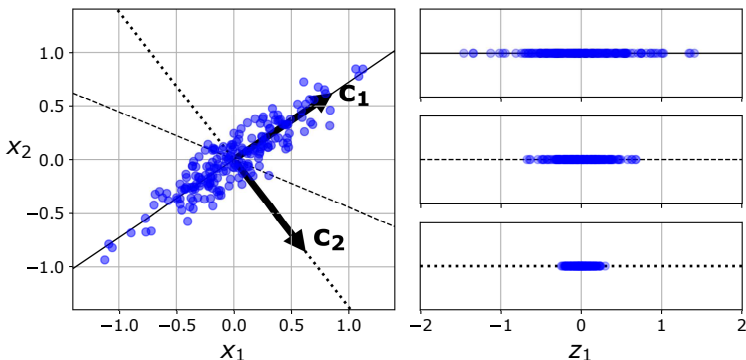
- 1 Identifica el hiperplano que se encuentra más cerca de los datos.
- 2 Luego proyecta los datos sobre él:



PCA: Preservando la varianza I

Antes de poder **proyectar** el conjunto de entrenamiento en un hiperplano de dimensiones inferiores, **primero debemos elegir el hiperplano correcto**.

- **Ejemplo:**

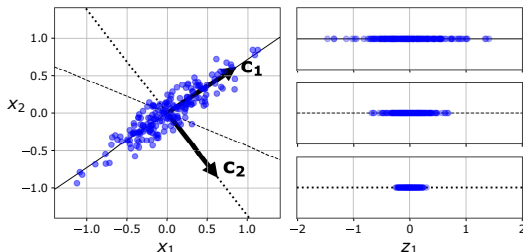


PCA: Preservando la varianza II

- En la **Figura de la izquierda** tenemos un conjunto de datos 2D, junto con tres ejes diferentes.
- En la **Figura de la derecha** está el resultado de la proyección de los datos en cada uno de estos ejes.
- Como podemos ver:
 - La proyección en la línea continua conserva la varianza máxima.
 - La proyección en la línea punteada conserva muy poca varianza.
 - La proyección en la línea discontinua conserva una cantidad intermedia de varianza.

Se selecciona el eje que **preserva la cantidad máxima de variación**, ya que perderá menos información. Es decir, el eje que **minimiza la distancia cuadrática media** entre el conjunto de datos original y su proyección sobre ese eje.

PCA: Componentes principales I

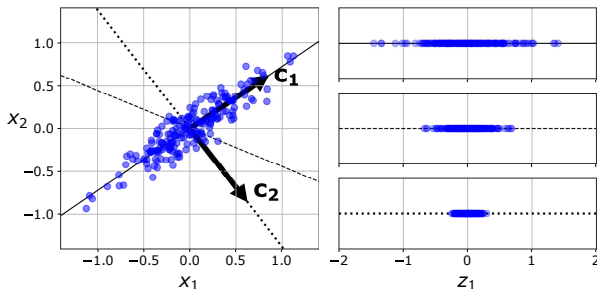


- PCA identifica el eje que representa la mayor cantidad de variación en el conjunto de entrenamiento (**línea continua**).
- También encuentra un segundo eje, **ortogonal al primero**, que representa la mayor cantidad de varianza restante (**línea punteada**).

PCA encontraría tantos ejes como el número de dimensiones en el conjunto de datos.

PCA: Componentes principales II

- El vector unitario que define el i -ésimo eje se llama i -ésimo componente principal (**PC**). En la Figura, la primera PC es c_1 y la segunda PC es c_2 .
- Las dos primeras PC están representadas por las flechas ortogonales en el plano, y la tercera PC sería ortogonal al plano (**apuntando hacia arriba o hacia abajo**).



PCA: Componentes principales III

- * **Tip:** La dirección de los componentes principales no es estable. Sin embargo, generalmente seguirán en los mismos ejes.

Entonces, ¿cómo encontrar los componentes principales de un conjunto de entrenamiento?

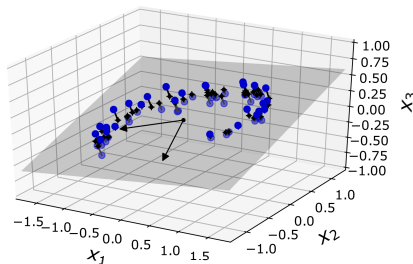
R/ Existe una técnica estándar de factorización de matriz llamada **Singular Value Decomposition (SVD)** que puede descomponer la matriz del conjunto de entrenamiento \mathbf{X} en la multiplicación matricial de tres matrices $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$, donde \mathbf{V} contiene todos los componentes principales que estamos buscando:

$$\mathbf{V} = \begin{pmatrix} | & | & \cdots & | \\ c_1 & c_2 & \cdots & c_n \\ | & | & & | \end{pmatrix}$$



PCA: Proyectando a d -dimensiones I

- Se reduce la dimensionalidad a d -dimensiones **proyectándolo en el hiperplano definido por los primero d componentes principales.** Garantizando que la proyección conservará la mayor variación posible.
- En la Figura, el conjunto de datos 3D se proyecta hacia abajo en el plano 2D definido por los dos primeros PCs, **preservando una gran parte de la varianza del conjunto de datos.**



PCA: Proyectando a d -dimensiones II

Para **proyectar el conjunto de entrenamiento en el hiperplano**, se calcula la multiplicación matricial de la matriz \mathbf{X} del conjunto de entrenamiento por la matriz \mathbf{W}_d , que contiene los primeros d componentes principales, como se muestra en la siguiente ecuación:

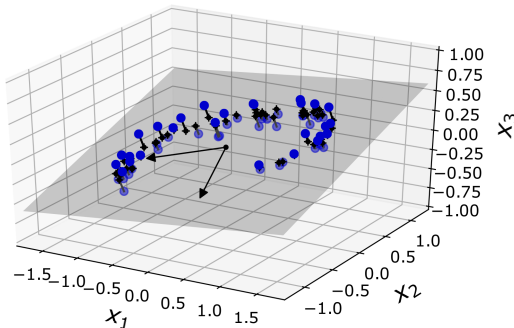
$$\mathbf{X}_{d-proj} = \mathbf{X}\mathbf{W}_d$$



PCA: Explained variance ratio

La **relación de varianza explicada** de cada componente principal, está disponible a través de la variable `explained_variance_ratio_`. Esta **indica la proporción de la varianza** del conjunto de datos que se encuentra a lo largo del eje de cada componente principal.

Ejemplo:





PCA: Elegir el número correcto de dimensiones

- **Generalmente es preferible elegir el número de dimensiones que se suman a una porción suficientemente grande de la varianza (por ejemplo, 95%).**



A menos, que el objetivo sea la visualización de datos; en ese caso, generalmente querrá reducir la dimensionalidad a 2 o 3.



- Obviamente, después de la **reducción de dimensionalidad**, el conjunto de entrenamiento ocupa mucho menos espacio.
- También es posible descomprimir el conjunto de datos reducido de nuevo a las dimensiones originales aplicando la transformación inversa de la proyección PCA. Por supuesto, esto no le devolverá los datos originales, ya que la proyección perdió un poco de información.

$$\mathbf{X}_{recovered} = \mathbf{X}_{d-proj} \mathbf{W}_d^T$$

La distancia cuadrática media entre los datos originales y los datos reconstruidos se denomina **error de reconstrucción**



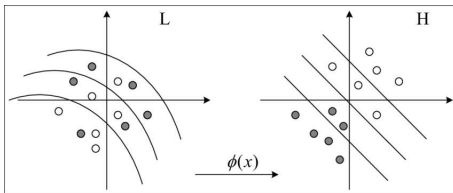
- **Randomized PCA:** Si configuramos el hiperparámetro `svd_solver` en "*randomized*", Scikit-Learn usa un algoritmo estocástico llamado PCA aleatorio que encuentra rápidamente una aproximación de los primeros d componentes principales.
- **Incremental PCA:** Un problema con las implementaciones anteriores de PCA es que requieren que todo el conjunto de entrenamiento se ajuste a la memoria para que el algoritmo se ejecute. PCA incremental (**IPCA**) divide el conjunto de entrenamiento en mini-lotes y alimenta un algoritmo de IPCA un mini-lote a la vez.



Kernel PCA

El *kernel trick*, mapea implícitamente las instancias en un espacio de muy alta dimensión (espacio de características), permitiendo la clasificación y regresión no lineal con SVMs.

- **TIP:** Recuerde que **un límite de decisión lineal** en el espacio de características de alta dimensión corresponde a un **límite de decisión no lineal complejo** en el espacio original.

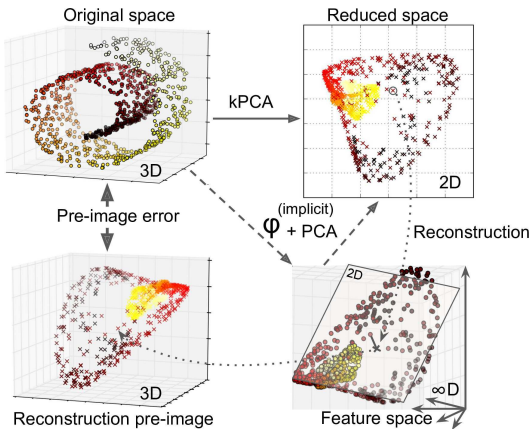


En PCA, se realizan proyecciones **no lineales complejas** para la reducción de la dimensionalidad. **Esto se llama Kernel PCA (kPCA)**



KPCA: selección de hiperparámetros I

Otro enfoque, esta vez **completamente sin supervisión**, es seleccionar el kernel y los hiperparámetros que producen **el error de reconstrucción más bajo**. Sin embargo, la reconstrucción no es tan fácil.



KPCA: selección de hiperparámetros II

- 1 Esto es **mapear los datos en un espacio de características de dimensiones infinitas** usando el feature map φ , y luego **proyectarlo usando PCA lineal**.
- 2 Al invertir el paso lineal de PCA para una instancia dada en el espacio reducido, **el punto reconstruido estaría en el espacio de características**, no en el original.
- 3 Dado que el espacio de características es de dimensión infinita, **no podemos calcular el punto reconstruido, ni el error de reconstrucción**.
- 4 Es posible encontrar un punto en el espacio original que se mapearía cerca del punto reconstruido. **Esto se llama reconstrucción pre-imagen**.

Se puede seleccionar el kernel y los hiperparámetros **que minimizan este error de reconstrucción pre-imagen**.



Contenido

- 1 Definición de reducción de dimensión
- 2 La maldición de la dimensionalidad
- 3 Enfoques principales para la DR
- 4 Análisis de componentes principales (PCA)
- 5 Locally Linear Embedding**
- 6 Otras técnicas de reducción de dimensión

Locally Linear Embedding (LLE)

- LLE es otra técnica de reducción de dimensionalidad **no lineal** (NLDR).
- LLE funciona midiendo primero **cómo cada instancia de entrenamiento se relaciona linealmente con sus vecinos más cercanos**.
- Luego busca una representación de baja dimensión del conjunto de entrenamiento donde estas **relaciones locales se conservan mejor**.
- Esto lo hace particularmente bueno para **desenrollar manifolds retorcidos**, especialmente cuando no hay demasiado ruido.

Ejemplo: Diseñemos un código que utilice la clase `LocallyLinearEmbedding` de Scikit-Learn para desenrollar el rollo suizo.



LLE: Cómo funciona? I

1. Para cada instancia de entrenamiento $\mathbf{x}^{(i)}$, el algoritmo identifica a sus k vecinos más cercanos. Luego intenta reconstruir $\mathbf{x}^{(i)}$ como una función lineal de estos vecinos. Específicamente, encuentra los pesos $w_{i,j}$ tal que la distancia al cuadrado entre $\mathbf{x}^{(i)}$ y $\sum_{j=1}^m w_{i,j} \mathbf{x}^j$ es lo más pequeña posible, suponiendo $w_{i,j} = 0$ si $\mathbf{x}^{(j)}$ no es uno de los k vecinos más cercanos de $\mathbf{x}^{(i)}$.

Por lo tanto, **el primer paso de LLE es un problema de optimización restringida**, donde \mathbf{W} es la matriz de peso que contiene todos los pesos $w_{i,j}$. La segunda restricción **simplemente normaliza los pesos** para cada instancia de entrenamiento $\mathbf{x}^{(i)}$.

$$\widehat{\mathbf{W}} = \arg \min_{\mathbf{W}} \sum_{i=1}^m \left(\mathbf{x}^{(i)} - \sum_{j=1}^m w_{i,j} \mathbf{x}^{(j)} \right)^2$$

sujeto a:

$$\begin{cases} w_{i,j} = 0 \\ \sum_{j=1}^m w_{i,j} = 1, \forall i = 1, 2, \dots, m \end{cases}$$

$w_{i,j} = 0$ si $\mathbf{x}^{(j)}$ no es uno de los k vecinos cercanos de $\mathbf{x}^{(i)}$. La matriz de peso \mathbf{W} **codifica las relaciones lineales locales entre las instancias de entrenamiento**.

- El segundo paso es **mapear las instancias de entrenamiento en un espacio d -dimensional** (donde $d < n$) mientras se **preservan estas relaciones locales** tanto como sea posible. Si $\mathbf{z}^{(i)}$ es la imagen de $\mathbf{x}^{(i)}$ en este espacio d -dimensional, entonces queremos que la distancia al cuadrado entre $\mathbf{z}^{(i)}$ y $\sum_{j=1}^m \widehat{w}_{i,j} \mathbf{z}^{(j)}$ **sea lo más pequeña posible**.

LLE: Cómo funciona? III

Esta idea nos lleva ahora a un problema de optimización **sin restricciones** descrito a continuación:

$$\hat{\mathbf{Z}} = \arg \min_{\mathbf{Z}} \sum_{i=1}^m \left(\mathbf{z}^{(i)} - \sum_{j=1}^m \hat{w}_{i,j} \mathbf{z}^{(j)} \right)^2$$

En lugar de mantener las instancias fijas y encontrar los pesos óptimos, estamos haciendo lo contrario: **mantener los pesos fijos y encontrar la posición óptima de las imágenes de las instancias en el espacio de baja dimensión**. Tengamos en cuenta que \mathbf{Z} es la matriz que contiene todo $\mathbf{z}^{(i)}$.

Desafortunadamente, el algoritmo LLE escala escasamente a conjuntos de datos muy grandes.

Contenido

- 1 Definición de reducción de dimensión
- 2 La maldición de la dimensionalidad
- 3 Enfoques principales para la DR
- 4 Análisis de componentes principales (PCA)
- 5 Locally Linear Embedding
- 6 Otras técnicas de reducción de dimensión

Otras técnicas de reducción de dimensión

- El **escalado multidimensional (MDS)** reduce la dimensionalidad al intentar preservar las distancias entre las instancias.
- **Isomap** crea un gráfico conectando cada instancia a sus vecinos más cercanos, luego reduce la dimensionalidad mientras intenta preservar las distancias geodésicas entre las instancias.
- **t-Distributed Stochastic Neighbor Embedding (t-SNE)** reduce la dimensionalidad al tratar de mantener instancias similares cercanas e instancias distintas lejos. Se utiliza principalmente para la visualización, en particular para visualizar grupos de instancias en espacios de alta dimensión.

Otras técnicas de reducción de dimensión



t-distributed stochastic neighbor embedding (t-SNE) I

- Dado un conjunto de N objetos de alta dimensión $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, t-SNE **calcula primero las probabilidades p_{ij} que son proporcionales a la similitud de los objetos \mathbf{x}_i y \mathbf{x}_j , de la siguiente manera:**

$$p_{j|i} = \frac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|\mathbf{x}_i - \mathbf{x}_k\|^2 / 2\sigma_i^2)} \quad (1)$$

- “La similitud del datapoint \mathbf{x}_j con el datapoint \mathbf{x}_i es la probabilidad condicional, $p_{j|i}$, que \mathbf{x}_i elegiría \mathbf{x}_j como vecino si los vecinos fueran seleccionados en proporción a su densidad de probabilidad bajo una Gaussiana centrada en \mathbf{x}_i ”.

$$p_{ij} = \frac{p_{j|i} + p_{i|j}}{2N} \quad (2)$$

t-distributed stochastic neighbor embedding (t-SNE) II

- Además, las probabilidades con $i = j$ se establecen en cero: $p_{ii} = 0$
- El ancho de banda del kernel Gaussiano se adapta a la densidad de los datos: se utilizan **valores más pequeños** de σ_i en partes **más densas** del espacio de datos.
- Dado que el kernel Gaussiano usa la distancia euclidiana $\|\mathbf{x}_i - \mathbf{x}_j\|$, se ve afectado por la **maldición de la dimensionalidad**, y en datos de alta dimensión **cuando las distancias pierden la capacidad de discriminar**, la p_{ij} se vuelve **demasiado similar**.

- t-SNE tiene como objetivo aprender un mapa d -dimensional $\{\mathbf{y}_1, \dots, \mathbf{y}_N\}$ (con $\mathbf{y}_i \in \mathbb{R}^d$) que **refleja las similitudes** p_{ij} lo mejor posible. Con este fin, mide las similitudes q_{ij} entre dos puntos en el mapa \mathbf{y}_i y \mathbf{y}_j , utilizando un enfoque muy similar. Específicamente, q_{ij} se define como:

$$q_{ij} = \frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_{k \neq i} (1 + \|\mathbf{y}_i - \mathbf{y}_k\|^2)^{-1}} \quad (3)$$

- Se usa una **distribución Student-t de cola pesada** para **medir similitudes entre puntos de baja dimensión para permitir que objetos diferentes se modelen muy separados en el mapa**. Tenga en cuenta que también en este caso configuramos $q_{ii} = 0$

- Las ubicaciones de los puntos \mathbf{y}_i en el mapa se determinan **minimizando la divergencia (no simétrica) Kullback–Leibler** de la distribución Q de la distribución P , es decir:

$$KL(P||Q) = \sum_{i \neq j} p_{ij} \log \frac{p_{ij}}{q_{ij}} \quad (4)$$

- La minimización de la divergencia *Kullback–Leibler* con respecto a los puntos \mathbf{y}_i se realiza usando **Gradiente Descendente**. El resultado de esta optimización es un mapa que refleja bien las similitudes entre las entradas de alta dimensión.