

DATOS MASIVOS II

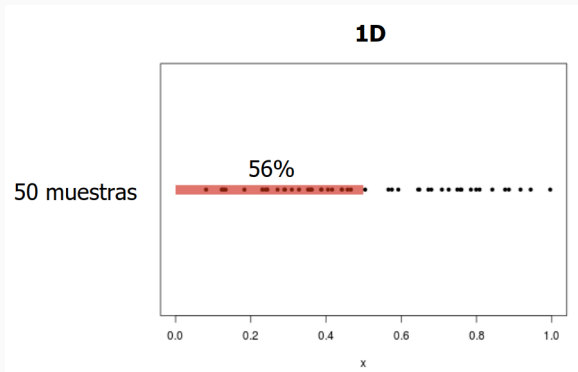
ANÁLISIS DE COMPONENTES PRINCIPALES

Blanca Vázquez-Gómez y Gibran Fuentes-Pineda

Septiembre 2020

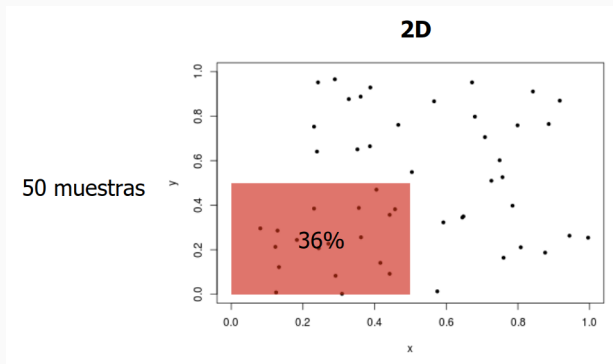
RECORDANDO LA MALDICIÓN DE LA DIMENSIONALIDAD

- Objetos cada vez más dispersos conforme aumenta el número de dimensiones



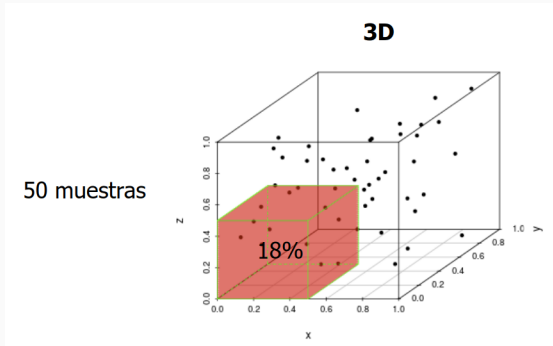
RECORDANDO LA MALDICIÓN DE LA DIMENSIONALIDAD

- Objetos cada vez más dispersos conforme aumenta el número de dimensiones



RECORDANDO LA MALDICIÓN DE LA DIMENSIONALIDAD

- Objetos cada vez más dispersos conforme aumenta el número de dimensiones



- Objetos cada vez más dispersos conforme aumenta el número de dimensiones
- **Ejercicio:** ¿Cuántas muestras necesitaría para cubrir un espacio de 1000 dimensiones con una precisión del 56 %?

LA HIPÓTESIS DE LA VARIEDAD

- Ejemplos pueden vivir en una variedad de muchas menores dimensiones que el espacio original

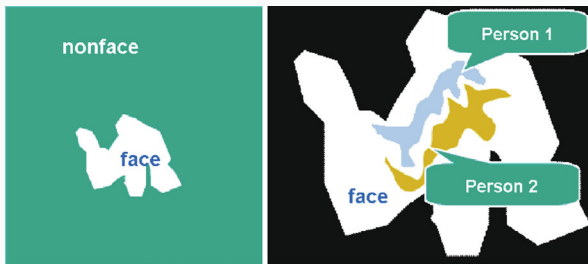
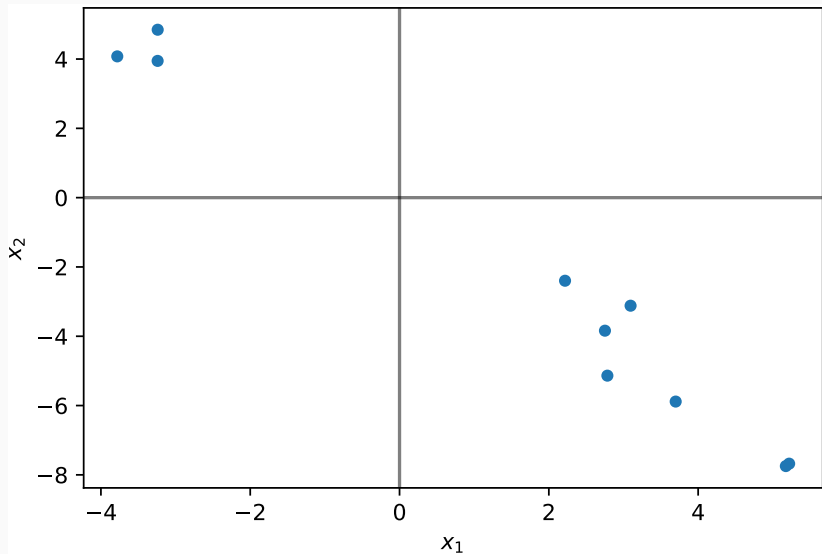


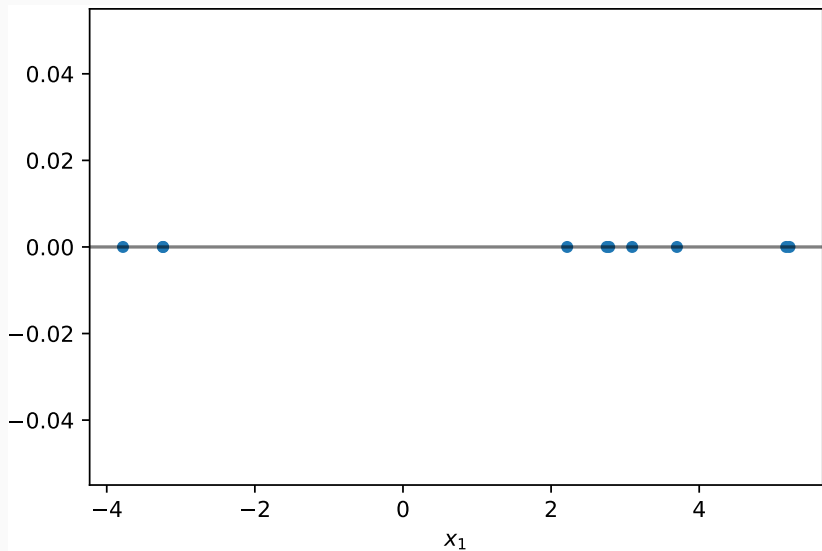
Imagen tomada de Li and Jain, 2005

- Proyección ortogonal de un conjunto de vectores
- Genera una nueva vista
- Aplicaciones
 - Visualización
 - Extracción de características
 - Reducción de dimensionalidad
 - Compresión

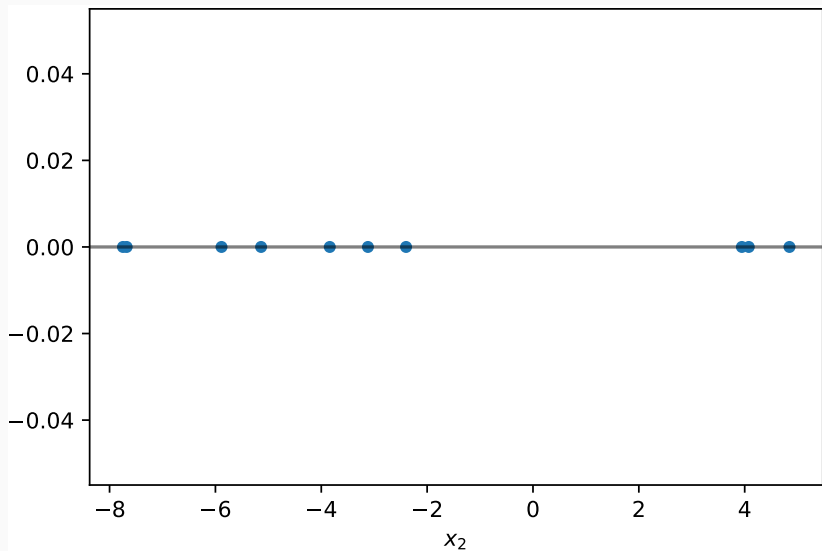
INTUICIÓN: DATOS EN 2D



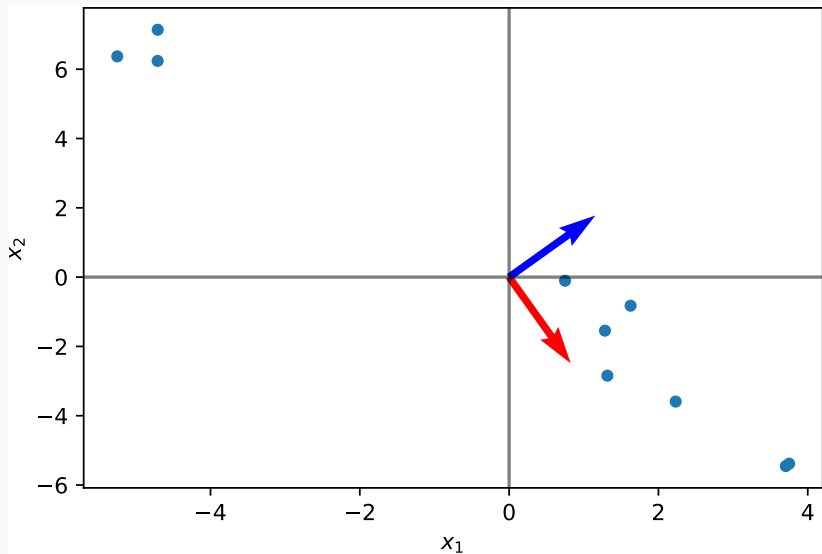
INTUICIÓN: DATOS VISTOS DESDE EL EJE x



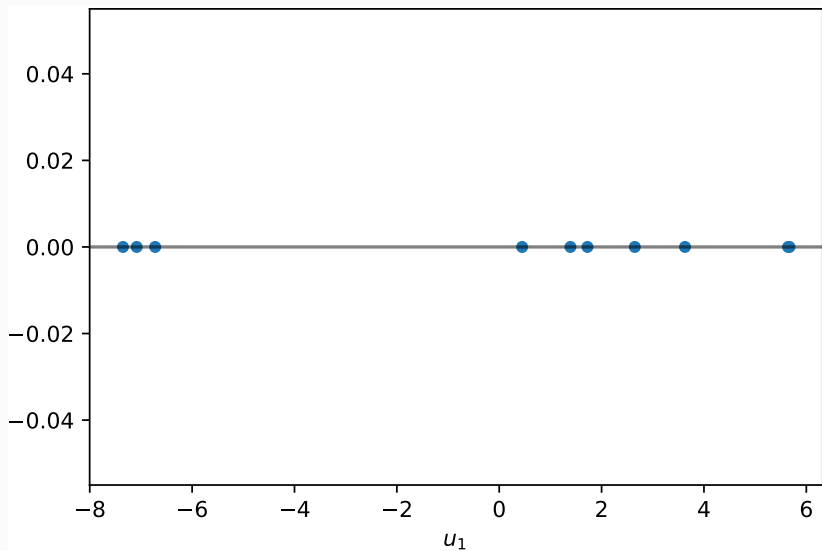
INTUICIÓN: DATOS VISTOS DESDE EL EJE y



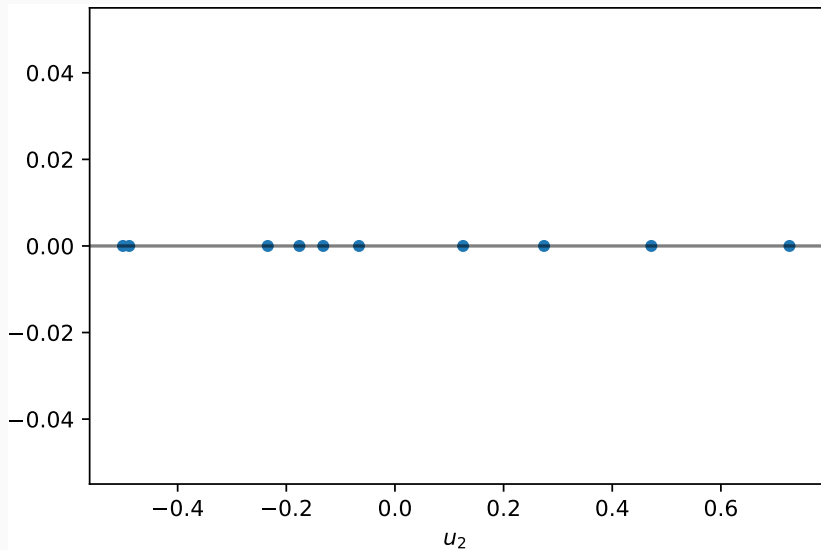
INTUICIÓN: NUEVOS EJES u_1 Y u_2



INTUICIÓN: DATOS PROYECTADOS SOBRE EL EJE u_1



INTUICIÓN: DATOS PROYECTADOS SOBRE EL EJE u_2



- Dado un conjunto de vectores $\mathbf{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}$ de d dimensiones, el primer componente principal es el vector \mathbf{u}_1 que maximice la varianza de los datos proyectados, donde \mathbf{u}_1 es un vector de d dimensiones

- La proyección de un vector sobre el componente principal está dado por

$$\hat{\mathbf{x}}^{(i)} = \mathbf{u}_1^\top \mathbf{x}^{(i)}$$

PCA POR MÁXIMA VARIANZA (1)

- La proyección de un vector sobre el componente principal está dado por

$$\hat{\mathbf{x}}^{(i)} = \mathbf{u}_1^\top \mathbf{x}^{(i)}$$

- La media de los datos proyectados es $\mathbf{u}_1^\top \bar{\mathbf{x}}$, donde

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)}$$

PCA POR MÁXIMA VARIANZA (1)

- La proyección de un vector sobre el componente principal está dado por

$$\hat{\mathbf{x}}^{(i)} = \mathbf{u}_1^\top \mathbf{x}^{(i)}$$

- La media de los datos proyectados es $\mathbf{u}_1^\top \bar{\mathbf{x}}$, donde

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)}$$

- La varianza es $\frac{1}{n} \sum_{i=1}^n [\mathbf{u}_1^\top \mathbf{x}^{(i)} - \mathbf{u}_1^\top \bar{\mathbf{x}}]^2 = \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1$, donde

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}^{(i)} - \bar{\mathbf{x}})(\mathbf{x}^{(i)} - \bar{\mathbf{x}})^\top$$

PCA POR MÁXIMA VARIANZA (2)

- Queremos encontrar el vector \mathbf{u}_1 que maximice la varianza de los datos proyectados $\mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1$, con la restricción que $\mathbf{u}_1^\top \mathbf{u}_1 = 1$

PCA POR MÁXIMA VARIANZA (2)

- Queremos encontrar el vector \mathbf{u}_1 que maximice la varianza de los datos proyectados $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$, con la restricción que $\mathbf{u}_1^T \mathbf{u}_1 = 1$
- Podemos formular esta restricción usando multiplicadores de Lagrange: $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda_1(1 - \mathbf{u}_1^T \mathbf{u}_1)$

PCA POR MÁXIMA VARIANZA (2)

- Queremos encontrar el vector \mathbf{u}_1 que maximice la varianza de los datos proyectados $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$, con la restricción que $\mathbf{u}_1^T \mathbf{u}_1 = 1$
- Podemos formular esta restricción usando multiplicadores de Lagrange: $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda_1(1 - \mathbf{u}_1^T \mathbf{u}_1)$
- Derivando e igualando a cero, tenemos

$$\mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1$$

PCA POR MÁXIMA VARIANZA (2)

- Queremos encontrar el vector \mathbf{u}_1 que maximice la varianza de los datos proyectados $\mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1$, con la restricción que $\mathbf{u}_1^\top \mathbf{u}_1 = 1$
- Podemos formular esta restricción usando multiplicadores de Lagrange: $\mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1 + \lambda_1(1 - \mathbf{u}_1^\top \mathbf{u}_1)$
- Derivando e igualando a cero, tenemos

$$\mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1$$

- Esto es, \mathbf{u}_1 es un vector propio de \mathbf{S} , donde $\lambda_1 = \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1$ es su valor propio que se corresponde con la varianza de los datos proyectados

- Para obtener el siguiente componente principal, buscamos el vector propio que maximice la varianza de los datos proyectados entre el conjunto de vectores ortogonales a los que ya han sido elegidos. Este proceso se realiza de forma incremental hasta obtener los m componentes principales.
- El conjunto de m componentes principales forman una base ortonormal de funciones.

- Gracias a que forman una base de funciones completa, podemos representar cualquier vector $\mathbf{x}^{(i)}$ como una combinación lineal de los componentes principales

$$\hat{\mathbf{x}}^{(i)} = \sum_{j=1}^m \alpha_{i,j} \mathbf{u}_j$$

donde $\alpha_{i,j} = \mathbf{x}^{(i)\top} \mathbf{u}_j$.

- Esto es

$$\hat{\mathbf{x}}^{(i)} = \sum_{j=1}^m (\mathbf{x}^{(i)\top} \mathbf{u}_j) \mathbf{u}_j$$

PCA POR VECTORES Y VALORES PROPIOS

- Busca subespacio de m dimensiones que maximiza varianza (o minimiza error) de los ejemplos
 - Definido por eigenvectores $\mathbf{u}_1, \dots, \mathbf{u}_m$ con eigenvalores más grandes $\lambda_1, \dots, \lambda_m$ de la matriz de covarianza

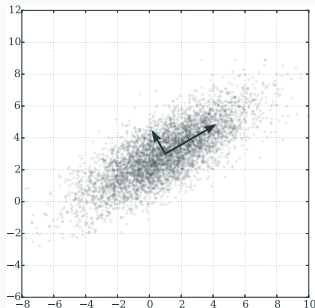
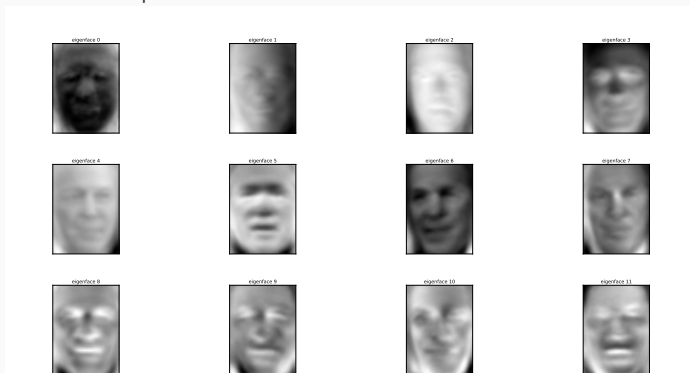


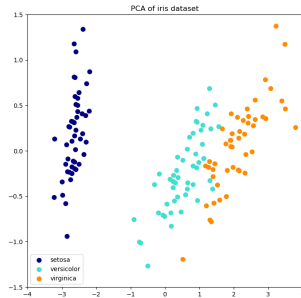
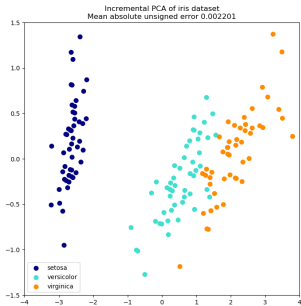
Figura tomada de Wikipedia (Principal Component Analysis)

PCA APLICADO A IMÁGENES DE ROSTROS

- Componentes principales se toman como base (**eigenfaces**)
- Nuevos rostros se proyectan en subespacio encontrado para ser comparados



PCA INCREMENTAL



Ejemplo de <http://scikit-learn.org>

ANÁLISIS DE FACTORES: VARIABLES CONTINUAS

- Variables latentes continuas $\mathbf{z} \in \mathbb{R}^K$, con a priori gaussiana

$$\mathbf{z} \sim \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$$

- Variables observadas continuas $\mathbf{x} \in \mathbb{R}^d$ con¹

$$\mathbf{x}|\mathbf{z} \sim \mathcal{N}(\mathbf{U}\mathbf{z} + \boldsymbol{\mu}, \boldsymbol{\Psi})$$

- Distribución sobre \mathbf{x} está dada por

$$P(\mathbf{x}) = \int P(\mathbf{x}|\mathbf{z})P(\mathbf{z})d\mathbf{z} = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{C})$$

donde $\mathbf{C} = \mathbf{U}\mathbf{U}^\top + \sigma^2\mathbf{I}$

¹Cuando $\boldsymbol{\Psi} = \sigma^2\mathbf{I}$, $\boldsymbol{\mu}_0 = \mathbf{0}$ y $\boldsymbol{\Sigma}_0 = \mathbf{I}$, se conoce como *análisis de componentes principales probabilista* (PPCA).

PROCESO GENERATIVO DE PPCA

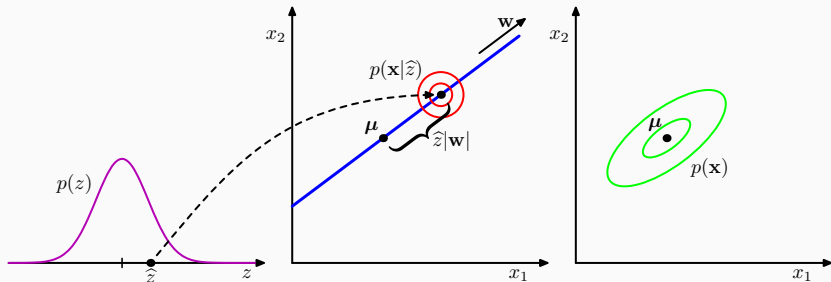


Imagen tomada de Bishop, PRML 2007

- Presuponiendo $\sigma^2 = 0$, se pueden encontrar parámetros de PCA por máxima verosimilitud usando el algoritmo EM
 1. Paso E: $\tilde{\mathbf{Z}} = (\mathbf{U}^\top \mathbf{U})^{-1} \mathbf{U}^\top \tilde{\mathbf{X}}$
 2. Paso M: $\mathbf{U} = \tilde{\mathbf{X}}^\top \tilde{\mathbf{Z}}^\top (\tilde{\mathbf{Z}} \tilde{\mathbf{Z}}^\top)^{-1}$donde $\tilde{\mathbf{X}} = \mathbf{X}^\top$