



# REDUCCIÓN DE DIMENSIONALIDAD

Dr. Jorge Hermosillo  
Laboratorio de Semántica Computacional



# VALORES PROPIOS Y VECTORES PROPIOS

---

# VALORES PROPIOS DE UNA MATRIZ

1. Definición: Un escalar  $\lambda$  es llamado un valor propio (*eigenvalue*) de la matriz  $A$   $n \times n$  si existe una solución no trivial de

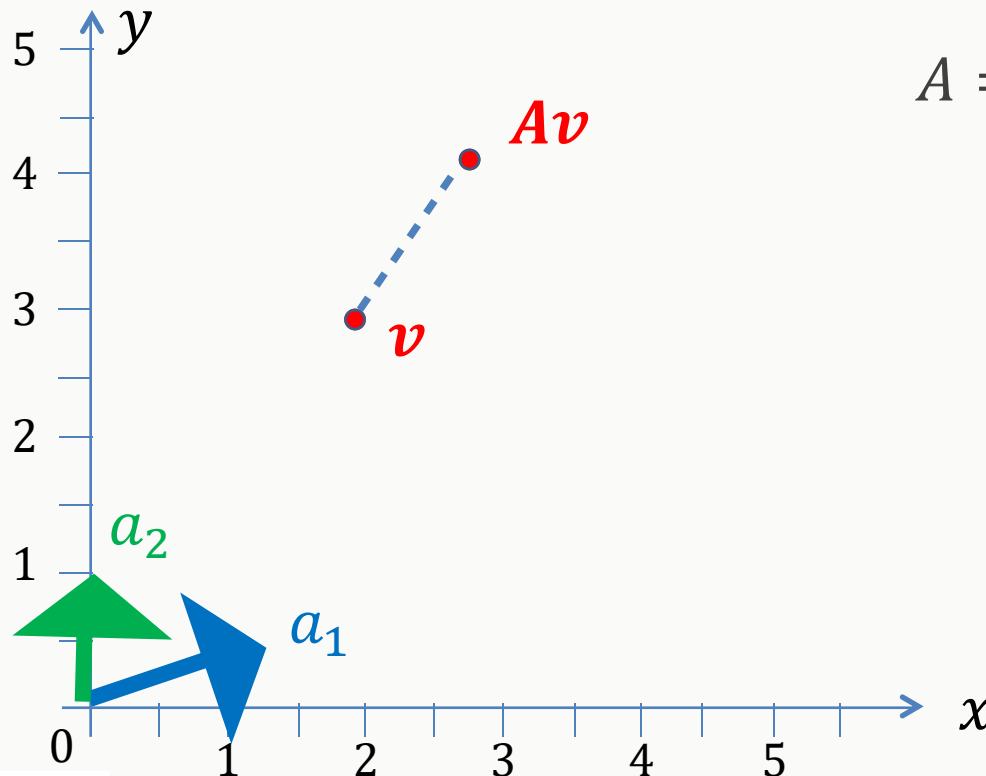
$$A\mathbf{x} = \lambda\mathbf{x}$$

tal que  $\mathbf{x}$  es llamado un ***vector propio*** (*eigenvector*) correspondiente al ***valor propio***  $\lambda$ .

¿A qué corresponde esto geométricamente?

# INTERPRETACIÓN GEOMÉTRICA DE VALORES Y VECTORES PROPIOS

- Sea  $v$  un vector y  $A$  una matriz con columnas  $a_1$  y  $a_2$  (mostradas como flechas). Si multiplicamos  $v$  por  $A$ , entonces  $A$  envía  $v$  a un nuevo vector  $Av$ .



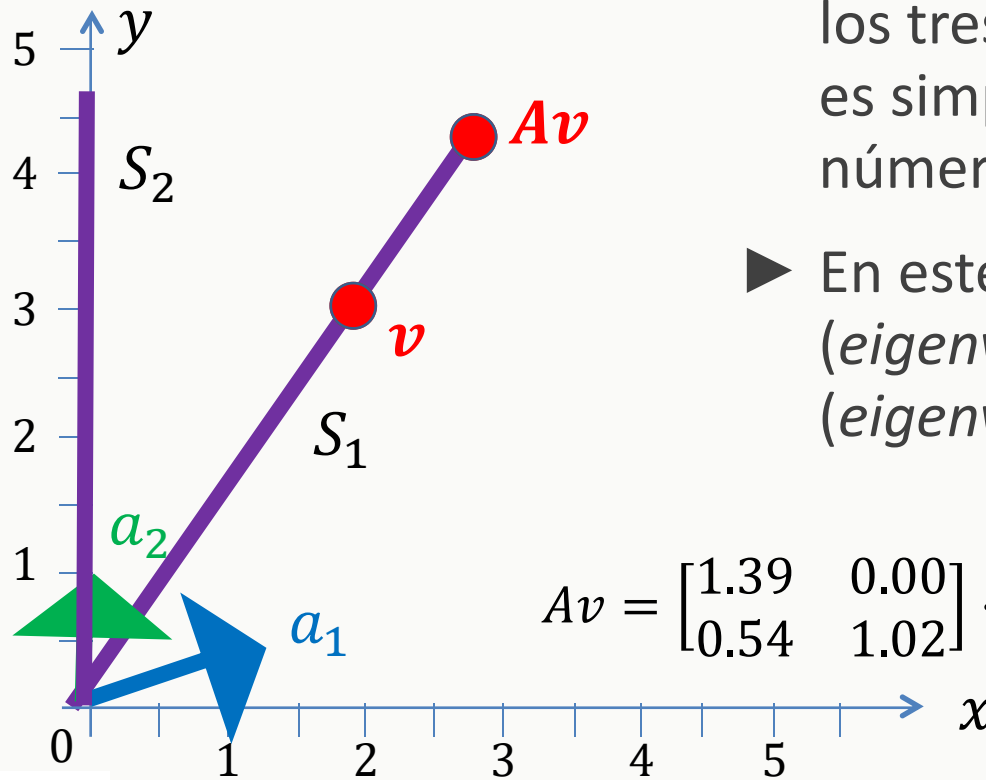
$$A = \begin{bmatrix} a_{1,x} & a_{2,x} \\ a_{1,y} & a_{2,y} \end{bmatrix} = \begin{bmatrix} 1.39 & 0.00 \\ 0.54 & 1.02 \end{bmatrix}$$

$$v = \begin{bmatrix} 2.00 \\ 2.91 \end{bmatrix}$$

$$Av = \begin{bmatrix} 2.78 \\ 4.04 \end{bmatrix}$$

<http://setosa.io/ev/eigenvectors-and-eigenvalues/>

# INTERPRETACIÓN GEOMÉTRICA DE VALORES Y VECTORES PROPIOS



- Si podemos dibujar una línea que pase por los tres puntos  $(0,0)$ ,  $v$  y  $Av$ , entonces  $Av$  es simplemente  $v$  multiplicado por un número  $\lambda$ ; es decir,  $Av = \lambda v$ .
- En este caso, llamamos  $\lambda$  un valor propio (*eigenvalue*) y  $v$  un vector propio (*eigenvector*).

$$Av = \begin{bmatrix} 1.39 & 0.00 \\ 0.54 & 1.02 \end{bmatrix} \cdot \begin{bmatrix} 2.00 \\ 2.91 \end{bmatrix} = \begin{bmatrix} 2.78 \\ 4.04 \end{bmatrix} = 1.39 \begin{bmatrix} 2.00 \\ 2.91 \end{bmatrix} \quad \begin{cases} \lambda_1 = 1.39 \\ \lambda_2 = 1.02 \end{cases}$$

# CONDICIÓN NECESARIA Y SUFICIENTE

- Si  $\lambda$  es un valor propio de la matriz  $A$   $n \times n$ , entonces existe una solución no trivial  $x$  de la ecuación

$$Ax = \lambda x$$

$$Ax - \lambda x = 0$$

$$Ax - \lambda Ix = 0$$

$$(A - \lambda I)x = 0$$

- Entonces  $x$  es un vector propio de  $A$  correspondiente al valor propio  $\lambda$  ssí  $x$  y  $\lambda$  satisfacen

$$(A - \lambda I)x = 0$$

# ECUACIÓN CARACTERÍSTICA

- ▶ La ecuación  $(A - \lambda I)x = 0$  tiene una solución no trivial ( $x \neq 0$ ) si y solo si la matriz  $A - \lambda I$  es no invertible (Teorema de Matriz Invertible).  
¿Por qué? Si la matriz es invertible entonces al multiplicar ambos lados por  $(A - \lambda I)^{-1}$  la única solución posible es  $x = 0$ .
- ▶ **Teorema:** El determinante de una matriz  $A$   $n \times n$  es 0 si y solo si la matriz  $A$  es no invertible.
- ▶ **Teorema:** Un escalar  $\lambda$  es un valor propio de una matriz  $A$   $n \times n$  si y solo si  $\lambda$  satisface la ecuación característica:

$$\det(A - \lambda I) = 0$$

# CARACTERÍSTICAS DE VALORES Y VECTORES PROPIOS

1. El espacio propio de una matriz  $A$   $n \times n$  correspondiente al valor propio  $\lambda$  de  $A$  es el conjunto de todos los vectores propios de  $A$  correspondientes a  $\lambda$ .
2. **Teorema:** si  $e_1, \dots, e_n$  son vectores propios que corresponden a distintos valores propios (i.e. multiplicidad = 1) de una matriz  $A$   $n \times n$ , entonces el conjunto  $\{e_1, \dots, e_n\}$  es linealmente independiente.
3. Si  $P = (e_1, \dots, e_n)$  es la matriz de vectores propios que corresponden a valores propios  $\lambda_1, \dots, \lambda_n$  entonces:

$$e_i^T e_j = 0, \lambda_i \neq \lambda_j \text{ (vectores ortogonales)}$$

$$P^T P = I$$

$$P^T = P^{-1}$$



# DIAGONALIZACIÓN DE UNA MATRIZ USANDO VECTORES PROPIOS

- Si  $A$  es una matriz  $n \times n$ , cuyos vectores propios son  $P = (e_1, \dots, e_n)$  y valores propios distintos son  $\{\lambda_1, \dots, \lambda_n\}$ .

Entonces:

$$A = P\Lambda P^T \Rightarrow P^T A P = \Lambda$$

donde  $\Lambda$  es la matriz diagonal:

$$\Lambda = \begin{bmatrix} \lambda_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_n \end{bmatrix}$$

# MATRIZ DE COVARIANZA Y VALORES PROPIOS

- ▶ Sea  $\Sigma$  la matriz de covarianza,  $\text{Tr}(\Sigma)$  –la traza de  $\Sigma$ – es la suma de los términos de la diagonal de  $\Sigma$
- ▶ Si  $\lambda_i, i = 1, 2, \dots, d$  son los valores propios de  $\Sigma$ , entonces

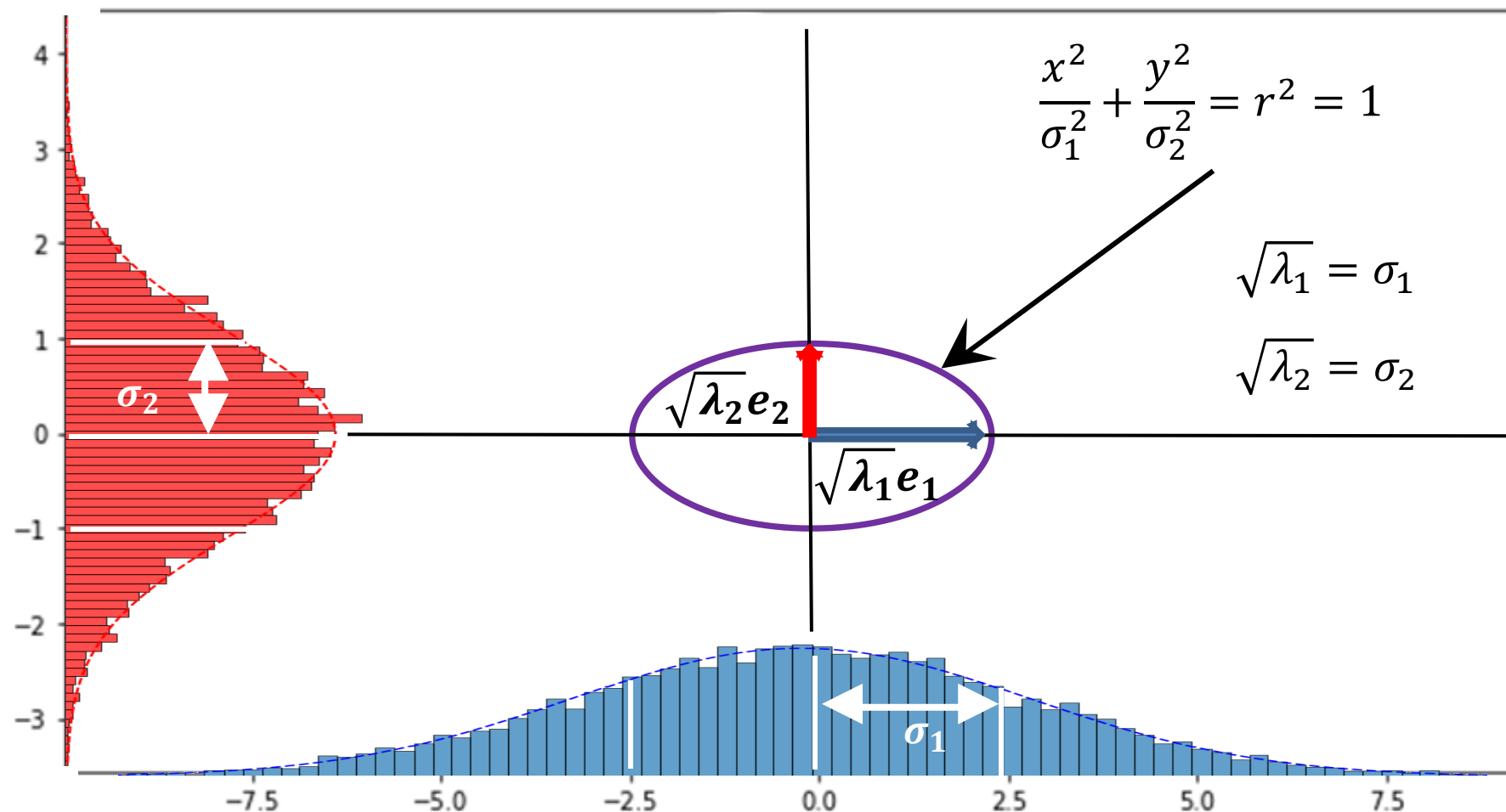
$$\sum_i \lambda_i = \text{Tr}(\Sigma)$$

i.e. los valores propios de  $\Sigma$  corresponden a las varianzas.

- ▶  $\det(\Sigma) = \prod_i \lambda_i$

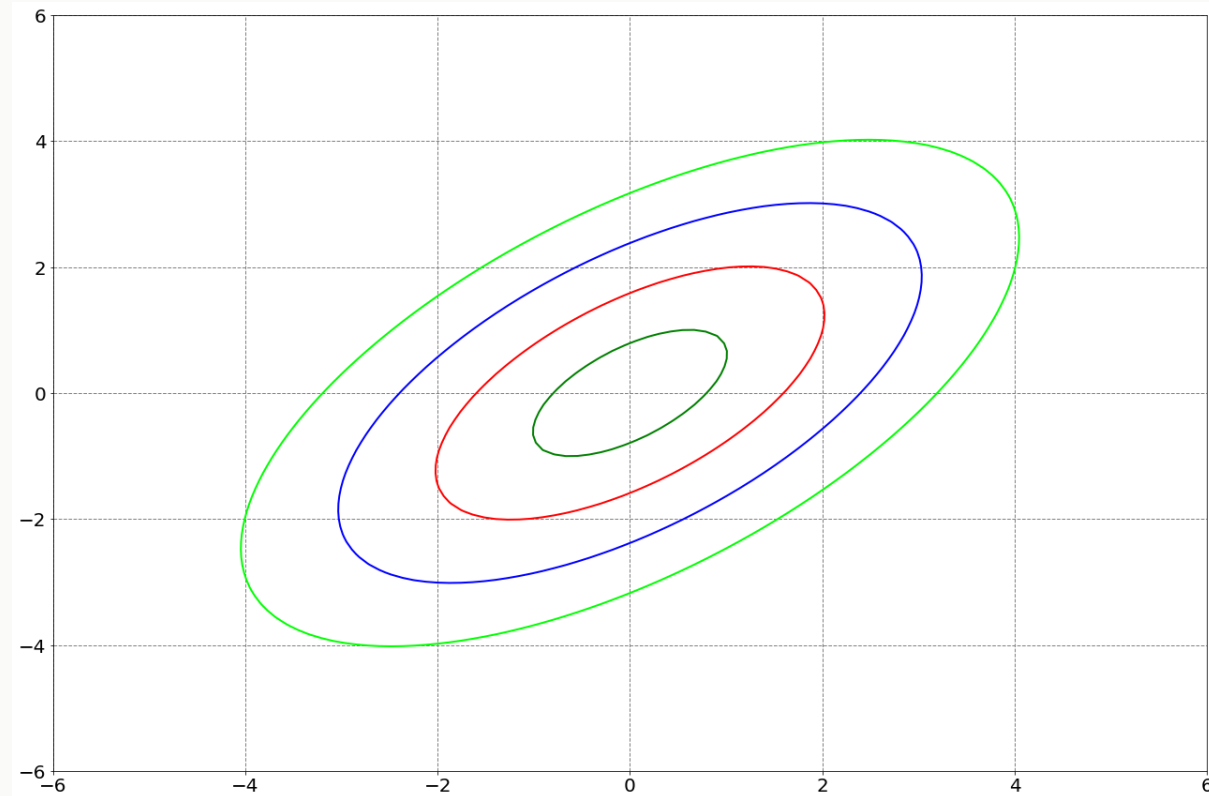
# INTERPRETACIÓN GEOMÉTRICA SIN CORRELACIÓN

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} = \begin{bmatrix} (2.5)^2 & 0 \\ 0 & 1 \end{bmatrix}$$



# INTERPRETACIÓN GEOMÉTRICA CON CORRELACIÓN

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} \sigma_1^2 & \rho_{12}\sigma_1\sigma_2 \\ \rho_{12}\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix} = \begin{bmatrix} 1 & 0.6 \\ 0.6 & 1 \end{bmatrix}$$



# DIAGONALIZACIÓN MATRIZ DE TRANSFORMACIÓN

► Diagonalización de una matriz (2D) mediante valores y vectores propios

- Supongamos la elipse es de la forma:  $\mathbf{x}^T \mathbf{M} \mathbf{x} = r^2$
- Pongamos  $\mathbf{M}$  en su forma diagonal usando sus vectores y valores propios:

$$\mathbf{M} = [\mathbf{e}_1 \ \mathbf{e}_2] \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} \begin{bmatrix} \mathbf{e}_1^T \\ \mathbf{e}_2^T \end{bmatrix}$$

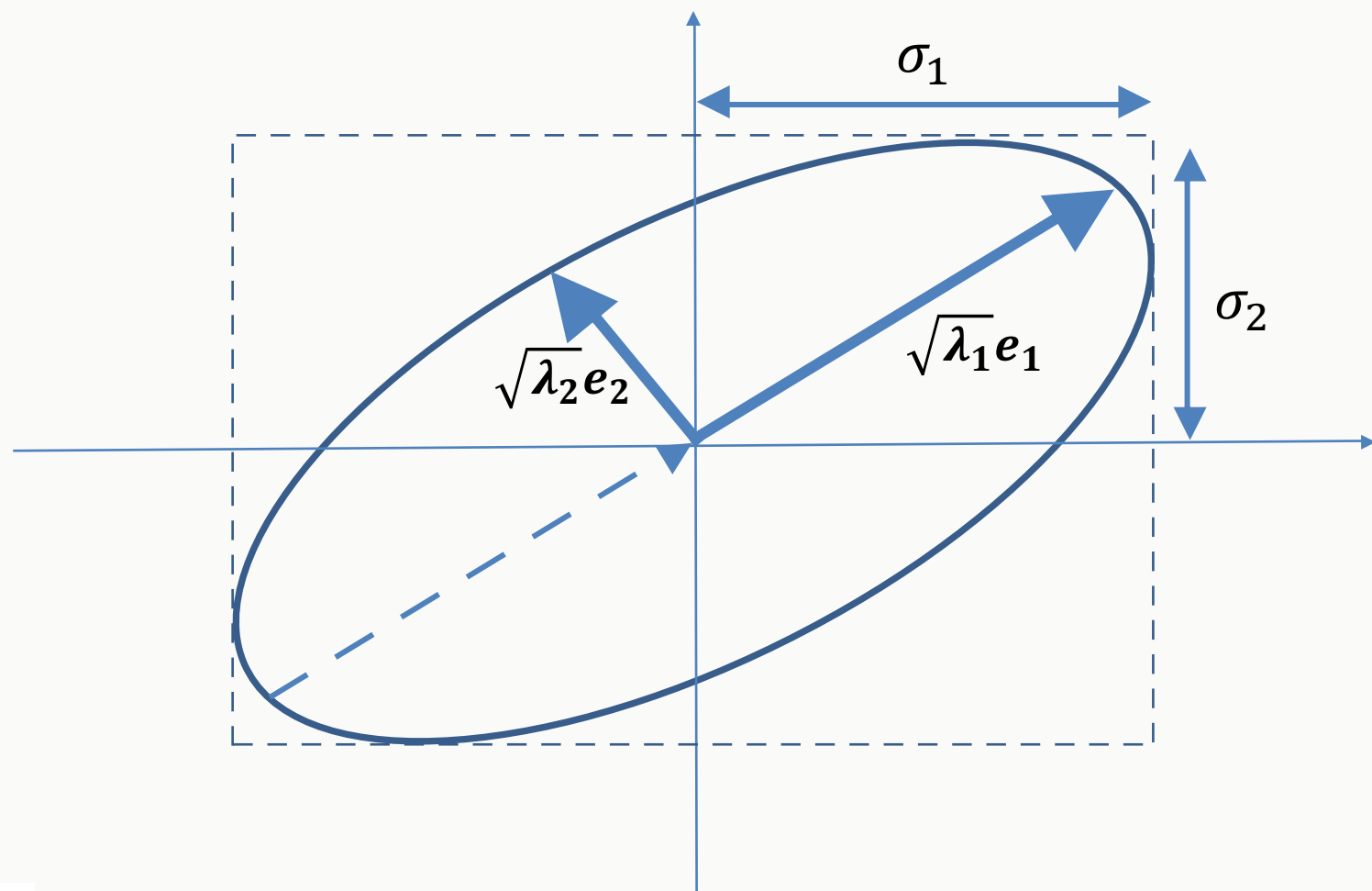
- La ecuación de la elipse es ahora:

$$\mathbf{x}^T \mathbf{M} \mathbf{x} = \mathbf{x}^T [\mathbf{e}_1 \ \mathbf{e}_2] \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} \begin{bmatrix} \mathbf{e}_1^T \\ \mathbf{e}_2^T \end{bmatrix} \mathbf{x} = r^2$$

$$\lambda_1 \mathbf{x}^T \mathbf{e}_1 \mathbf{e}_1^T \mathbf{x} + \lambda_2 \mathbf{x}^T \mathbf{e}_2 \mathbf{e}_2^T \mathbf{x} = r^2$$

$$\frac{(\mathbf{e}_1^T \mathbf{x})^2}{(1/\sqrt{\lambda_1})^2} + \frac{(\mathbf{e}_2^T \mathbf{x})^2}{(1/\sqrt{\lambda_2})^2} = r^2 \text{ donde } \lambda_1 \text{ es el valor propio más pequeño}$$

# INTERPRETACIÓN GEOMÉTRICA CON CORRELACIÓN



# PRINCIPAL COMPONENT ANALYSIS (PCA)

---

# REDUCCIÓN DE DIMENSIONALIDAD

- ▶ La reducción de dimensionalidad es muy importante en aplicaciones de ciencia de datos, bioinformática, recuperación de información, machine learning y reconocimiento de patrones.
- ▶ Existen dos enfoques: supervisado (e.g. Mixture Discriminant Analysis (MDA); Linear Discriminant Analysis (LDA)), y no supervisado (e.g. Independent Component Analysis (ICA); Principal Component Analysis (PCA))



# ¿PARA QUÉ USAR PCA?

► PCA se usa para:

- *Encontrar relaciones entre observaciones.*
- *Extraer la información más importante de los datos.*
- *Detectar y remover de valores atípicos (“ouliers”).*
- *Reducir la dimensionalidad para fines de visualización.*
- *Obtener features para fines de clasificación.*

## Objetivo de PCA

- ▶ Encontrar un espacio de dimensión reducida que se usa para transformar datos de un espacio de dimensión más alta a un espacio de dimensión más pequeña.
- ▶ Encontrar factores comunes, los llamados componentes principales, en forma de combinaciones lineales de las variables bajo investigación, y ordenarlas.
- ▶ **Identificar la base (vectorial) más significativa para volver a expresar un conjunto de datos.**

# EJEMPLO

- El objetivo es obtener una simple ecuación de  $x$ . Pero...
- Los datos están referenciados a los ejes de coordenadas de cada cámara y los desconocemos.
- Hay ruido en las mediciones
- No sabemos cuál de las señales refleja mejor las propiedades del sistema.
- PCA será útil para conocer  $x$ .

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

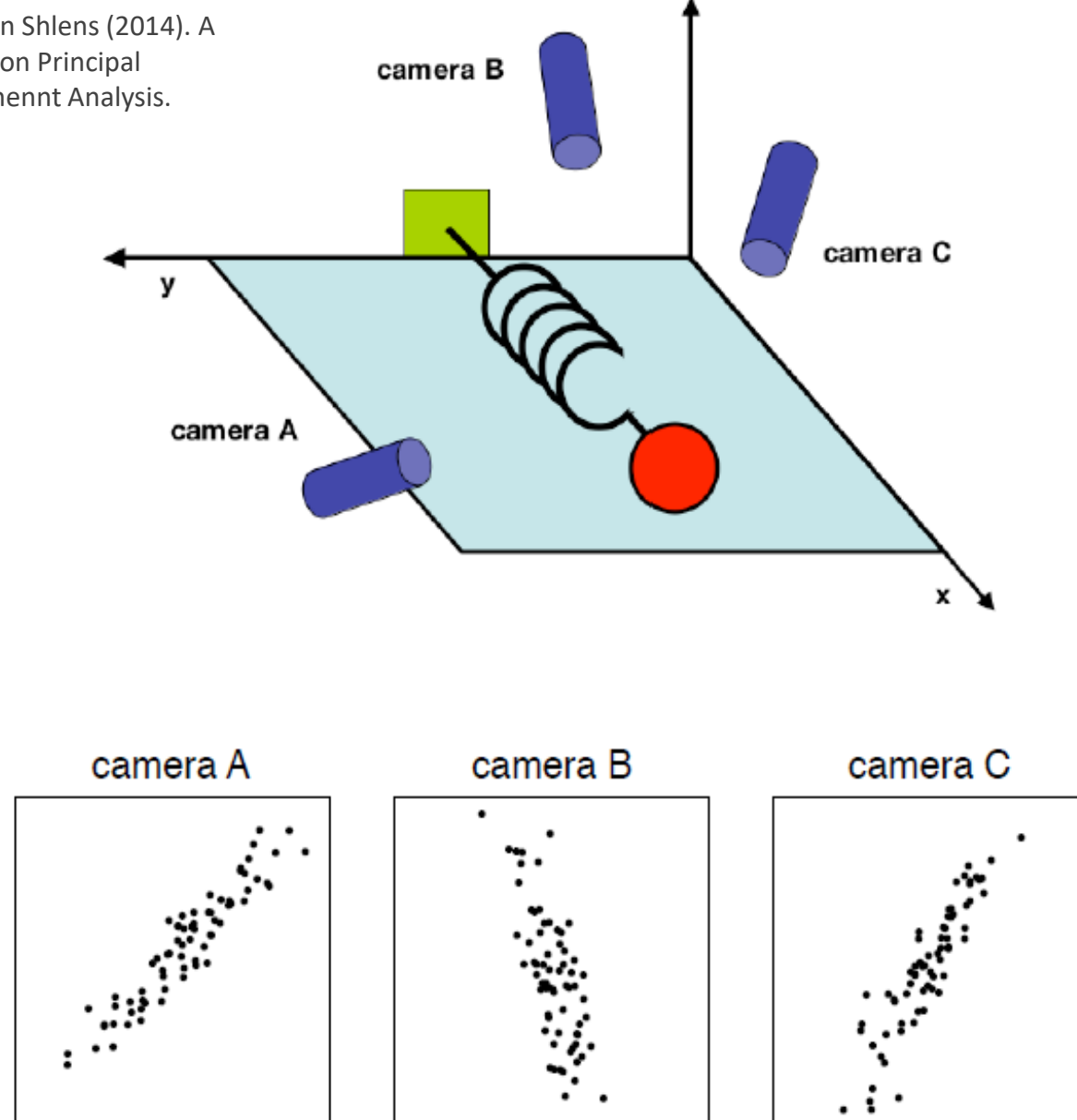
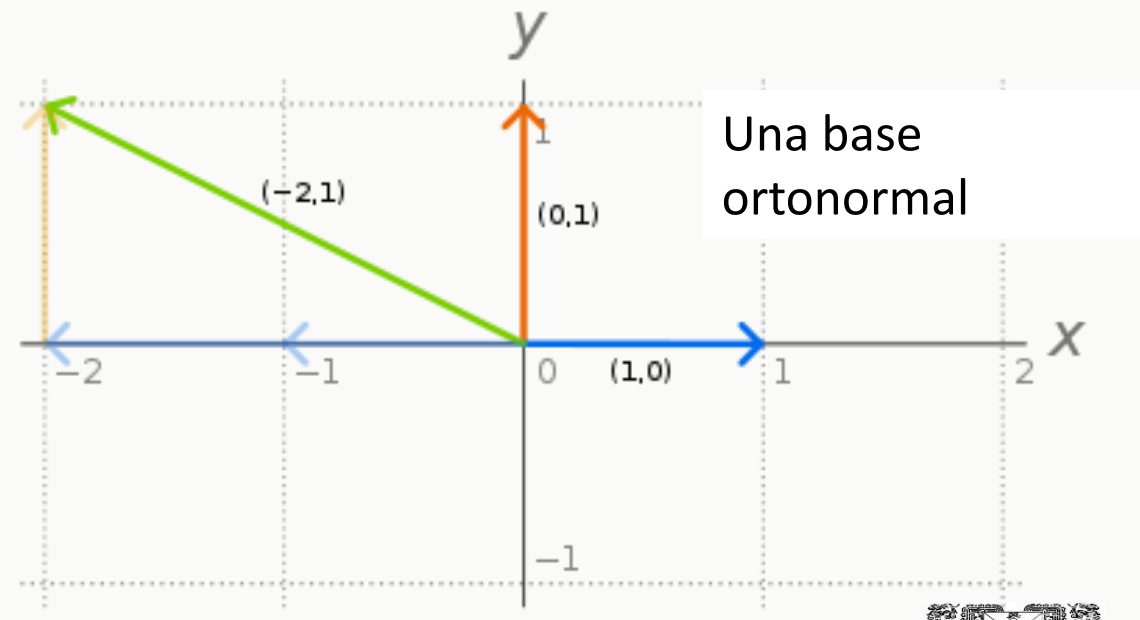
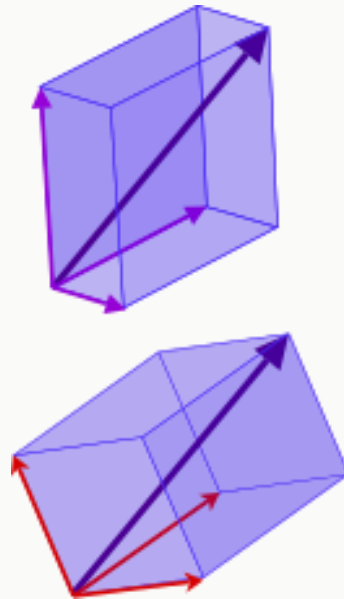


FIG. 1 A toy example. The position of a ball attached to an oscillating spring is recorded using three cameras A, B and C. The position of the ball tracked by each camera is depicted in each panel below.

## PCA → CAMBIO DE BASE

- La idea es que PCA nos ayudará a cambiar de base y que en esta nueva base filtraremos el ruido y mostraremos la estructura interna; para el ejemplo: “la dinámica ocurre sobre el eje  $x$ ”, es decir, la dimensión importante es el eje  $x$ .

Dos bases para un mismo vector



# PUNTO DE PARTIDA: UNA BASE “INGENUA” PARA LOS DATOS

- Supongamos que nuestros datos:

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$$

- Están referenciados a una base:

$$B = \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix} = \begin{bmatrix} 1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1 \end{bmatrix} = \mathbf{I}$$

- Toda observación  $X$  sería una combinación lineal de  $(b_1, \dots, b_m)$ .
- Esta base refleja el método de medición de nuestros datos.

# CAMBIO DE BASE

► La pregunta es ahora:

*¿existe otra base, que sea una combinación lineal de la base original, que re-exprese mejor nuestro conjunto de datos?*

► PCA supone linealidad: está limitado a re-expresar los datos como una **combinación lineal** de sus vectores base.

## PRINCIPIO BÁSICO

- ▶ Sea  $\mathbf{X}$  la matriz  $m \times n$  de datos originales, donde cada columna es un vector dimensión  $m$ ; por lo tanto,  $n$  es el número de observaciones.
- ▶ Sea  $\mathbf{Y}$  otra matriz  $m \times n$  relacionada a  $\mathbf{X}$  por la transformación  $\mathbf{P}$ .
- ▶  $\mathbf{X}$  es la matriz de datos originales, y  $\mathbf{Y}$  es la nueva representación del conjunto de datos:

$$\mathbf{PX} = \mathbf{Y} \quad (1)$$

# INTERPRETACIONES DE PCA

- ▶ La ecuación (1) representa un cambio de base y como tal puede tener muchas interpretaciones:
  - $\mathbf{P}$  es una matriz que transforma  $\mathbf{X}$  en  $\mathbf{Y}$ .
  - Geométricamente,  $\mathbf{P}$  es una rotación y escalamiento que, nuevamente, transforma  $\mathbf{X}$  en  $\mathbf{Y}$ .
  - Los renglones de  $\mathbf{P}$ ,  $\{\mathbf{p}_1, \dots, \mathbf{p}_m\}$ , son un conjunto de nuevos vectores base para expresar las columnas de  $\mathbf{X}$ .



LOS RENGLONES DE  $\mathbf{P}$  SON UN CONJUNTO DE NUEVOS VECTORES BASE PARA EXPRESAR LAS COLUMNAS DE  $\mathbf{X}$

$$\mathbf{P}\mathbf{X} = \begin{bmatrix} \mathbf{p}_1 \\ \vdots \\ \mathbf{p}_m \end{bmatrix} [\mathbf{x}_1 \quad \cdots \quad \mathbf{x}_n]$$

$$\mathbf{Y} = \begin{bmatrix} \mathbf{p}_1 \cdot \mathbf{x}_1 & \cdots & \mathbf{p}_1 \cdot \mathbf{x}_n \\ \vdots & \ddots & \vdots \\ \mathbf{p}_m \cdot \mathbf{x}_1 & \cdots & \mathbf{p}_m \cdot \mathbf{x}_n \end{bmatrix}$$

Cada coeficiente  $y_i$  de  $\mathbf{Y}$  es un producto punto de  $\mathbf{x}_i$  con el renglón correspondiente en  $\mathbf{P}$ .

En otras palabras, el coeficiente  $j$  de  $y_i$  es una proyección sobre el renglón  $j$  de  $\mathbf{P}$ .

## LO QUE QUEDA POR RESOLVER

- ▶ El problema se reduce a encontrar el apropiado cambio de base.
- ▶ Los vectores fila  $\{\mathbf{p}_1, \dots, \mathbf{p}_m\}$  en esta transformación serán los componentes principales de  $\mathbf{X}$ .
  - *¿Cuál es la mejor manera de volver a expresar  $\mathbf{X}$ ?*
  - *¿Cuál es una buena opción de base  $\mathbf{P}$ ?*
- ▶ Estas preguntas se responden preguntándonos qué características nos gustaría que  $\mathbf{Y}$  exhibiera.

# RUIDO Y ROTACIÓN

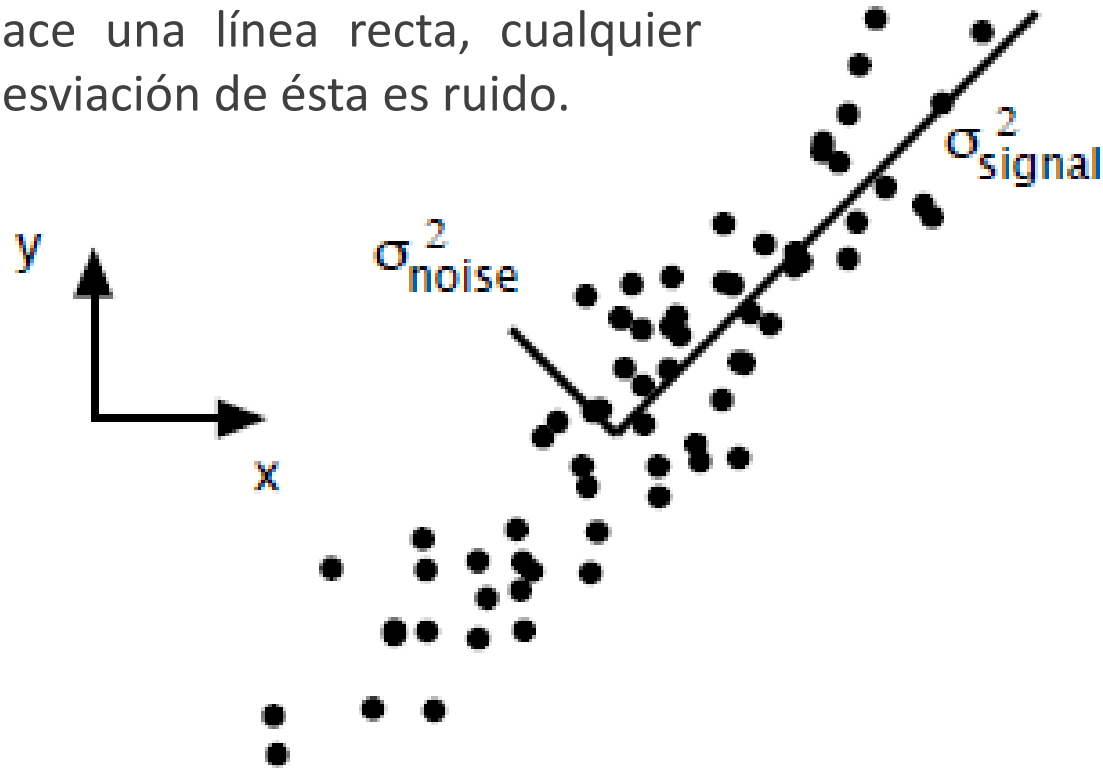
- En Teoría de la Información una medida común de la calidad de una señal es la *Relación Señal a Ruido* (SNR por sus siglas en inglés).

$$SNR = \frac{\sigma_{señal}^2}{\sigma_{ruido}^2}$$

- Un alto valor de SNR ( $\gg 1$ ) indica una medición de precisión.

# RUIDO Y ROTACIÓN: EL EJEMPLO DE LA CÁMARA

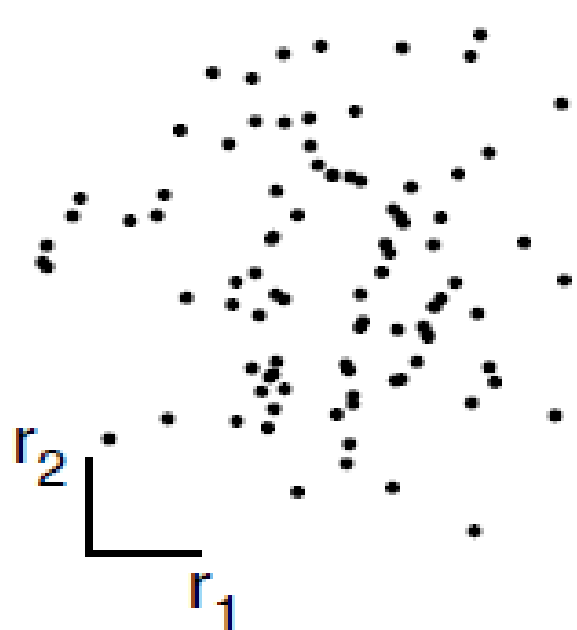
Si suponemos que el resorte hace una línea recta, cualquier desviación de ésta es ruido.



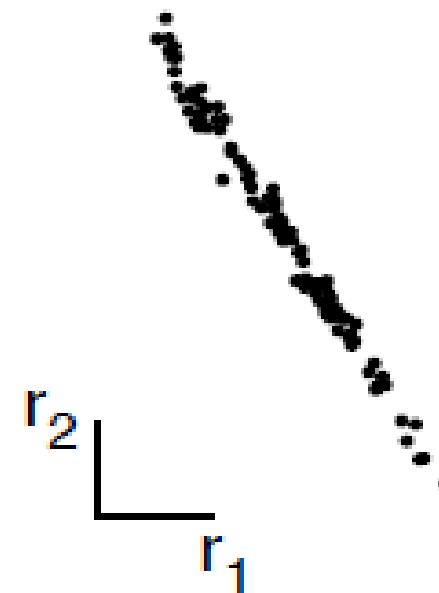
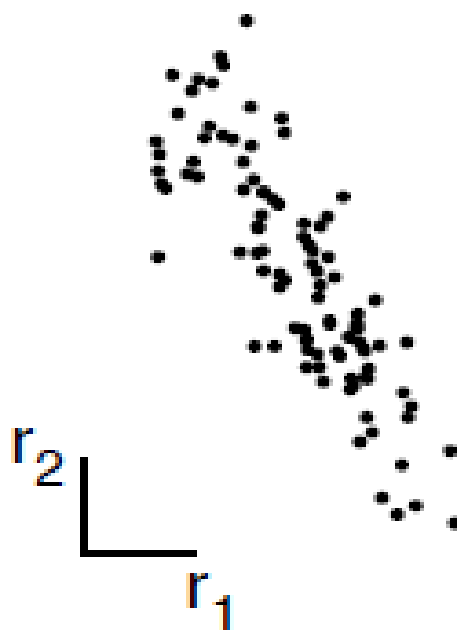
Jonathon Shlens (2014). A tutorial on Principal Component Analysis.

- Nota que la dirección de la varianza más grande no está en la dirección de la base del registro  $(x_A, y_A)$  sino más bien sobre la línea de mejor ajuste; i.e. la información importante está ahí
- → rotar la base original (ingenua) hacia la línea de mejor ajuste revelaría la dirección del movimiento del resorte en 2-D.

# REDUNDANCIA



low redundancy



high redundancy

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

# MATRIZ DE COVARIANZA

- La matriz de covarianza permite generalizar la noción de redundancia a un número arbitrario de dimensiones.

$$\mathbf{C}_X = \frac{1}{N} \mathbf{X} \mathbf{X}^T$$

- *En los términos de la diagonal, como supuesto, valores altos corresponden a estructura interesante.*
- *En los términos fuera de la diagonal, magnitudes grandes corresponden a alta redundancia.*

# DIAGONALIZACIÓN DE LA MATRIZ DE COVARIANZA

- De acuerdo con lo anterior, para nuestra matriz  $\mathbf{C}_Y$  deseamos: (1) minimizar la redundancia, medida por la magnitud de la covarianza, y (2) maximizar la señal, medida por la varianza

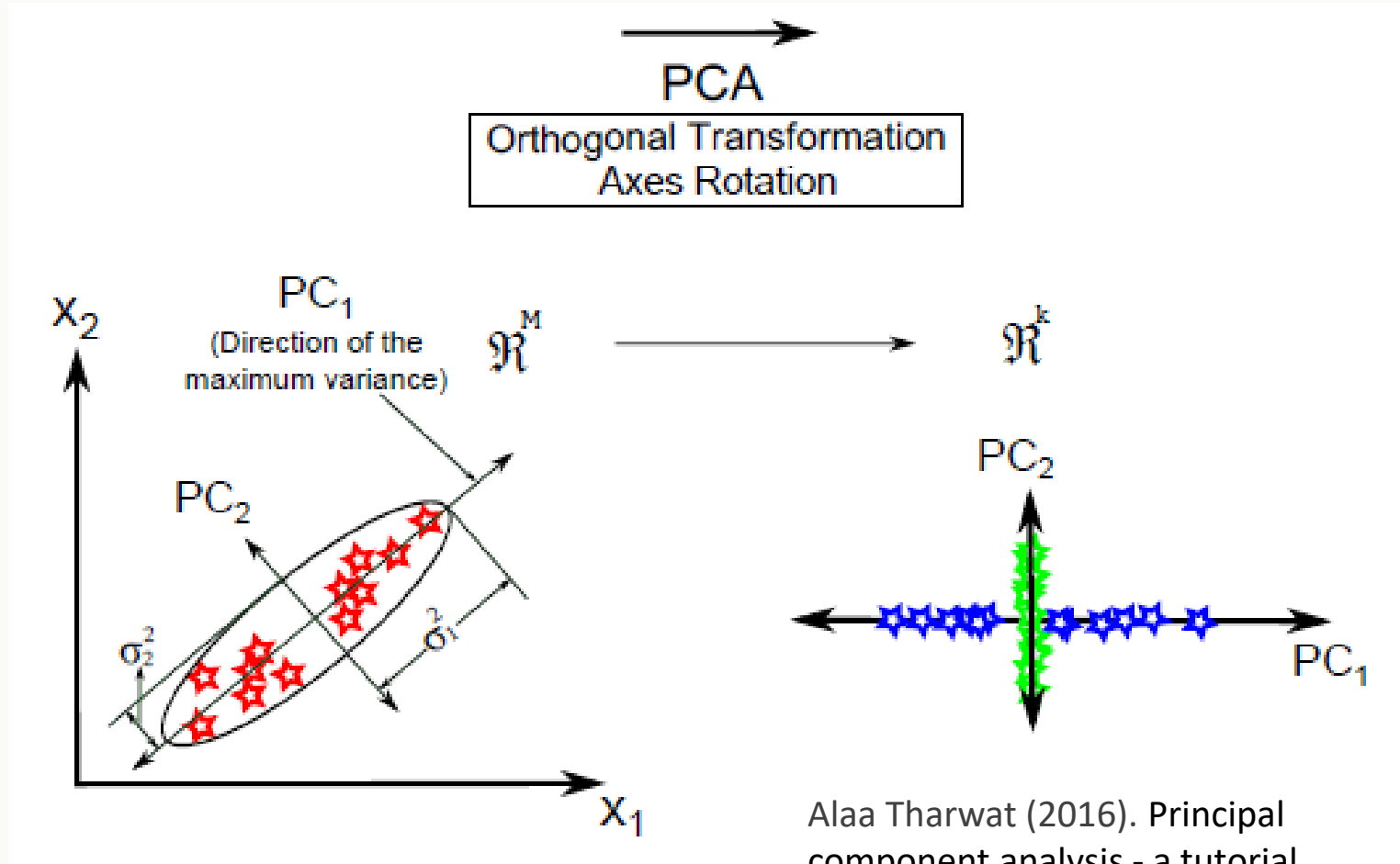
*¡¡Hay que diagonalizar  $\mathbf{C}_Y$ !!*

# SUPUESTOS DE PCA

- ▶ *Linealidad*: esto enmarca el problema como un cambio de base.
- ▶ *Varianza grande tiene estructura importante*: esto abarca la suposición de que los datos tienen una alta SNR. Por lo tanto, las componentes principales con mayores asociadas varianzas representan estructura interesante, mientras que, aquellas con poca varianza representan ruido.
- ▶ *Las componentes principales son ortogonales*: una suposición que facilita el cálculo mediante técnicas de álgebra lineal.



# SOLUCIÓN UTILIZANDO VECTORES PROPIOS



OBJETIVO:

Encuentra una matriz ortonormal  $\mathbf{P}$  en  $\mathbf{Y} = \mathbf{P}\mathbf{X}$  tal que  $\mathbf{C}_Y = \frac{1}{n}\mathbf{Y}\mathbf{Y}^T$  es una matriz diagonal.

Las filas de  $\mathbf{P}$  son las componentes principales de  $\mathbf{X}$ .

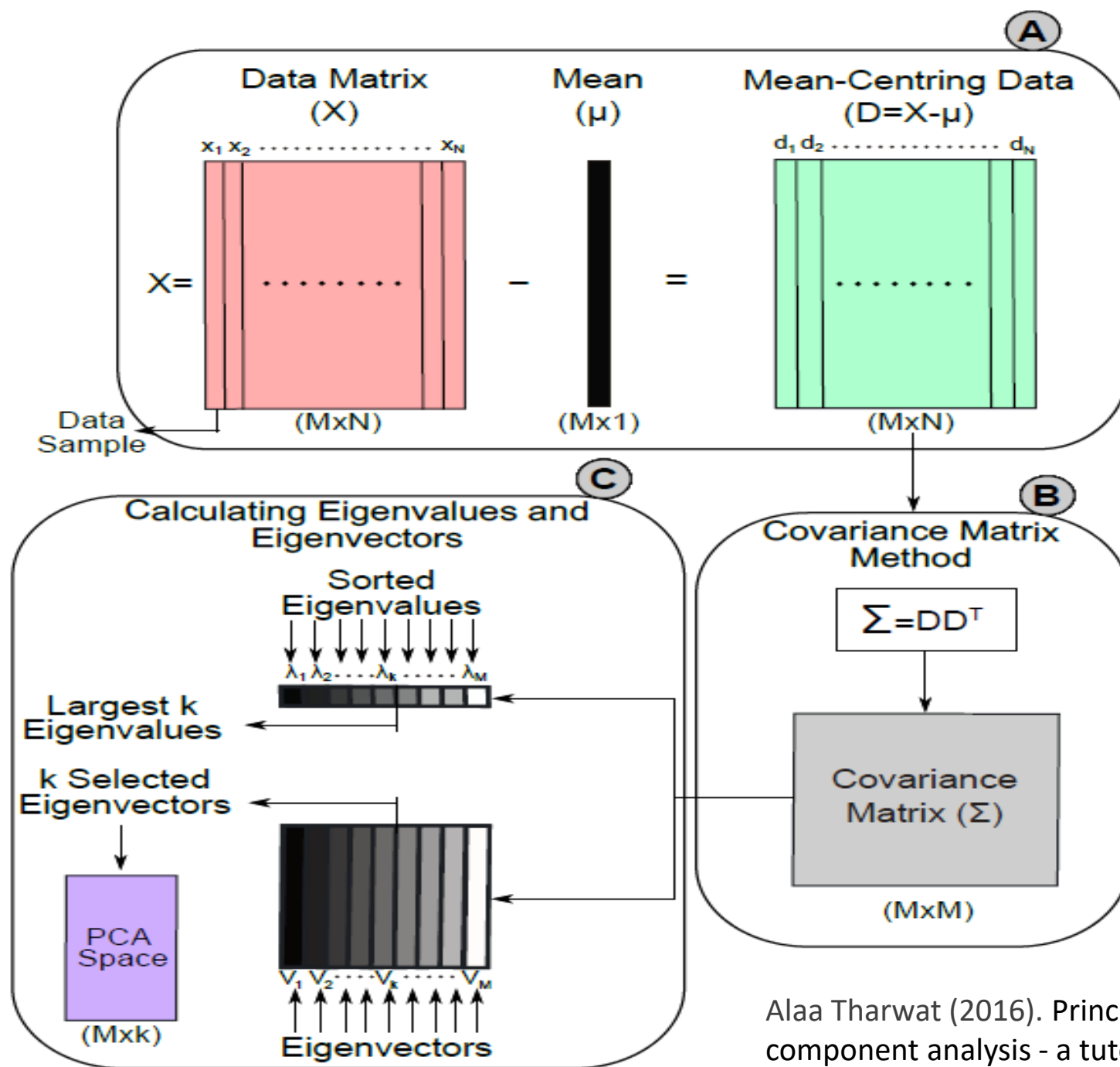
## EXPRESIÓN DE $\mathbf{C}_Y$ EN FUNCIÓN DE $\mathbf{C}_X$

$$\begin{aligned}\mathbf{C}_Y &= \frac{1}{n} \mathbf{Y} \mathbf{Y}^T \\ &= \frac{1}{n} (\mathbf{P} \mathbf{X}) (\mathbf{P} \mathbf{X})^T \\ &= \frac{1}{n} \mathbf{P} \mathbf{X} \mathbf{X}^T \mathbf{P}^T \\ &= \mathbf{P} \frac{1}{n} \mathbf{X} \mathbf{X}^T \mathbf{P}^T \\ \mathbf{C}_Y &= \mathbf{P} \mathbf{C}_X \mathbf{P}^T\end{aligned}$$

## DESCOMPOSICIÓN DE $\mathbf{C}_X$ EN VECTORES PROPIOS $\mathbf{E}$ (ASUMIMOS QUE $\mathbf{P} \equiv \mathbf{E}$ )

$$\begin{aligned}\mathbf{C}_Y &= \mathbf{P}\mathbf{C}_X\mathbf{P}^T \\ &= \mathbf{P}(\mathbf{E}^T\mathbf{D}\mathbf{E})\mathbf{P}^T \\ &= \mathbf{P}(\mathbf{P}^T\mathbf{D}\mathbf{P})\mathbf{P}^T \\ &= \mathbf{P}\mathbf{P}^{-1}\mathbf{D}\mathbf{P}\mathbf{P}^{-1} \\ \mathbf{C}_Y &= \mathbf{D}\end{aligned}$$

- Las componentes principales de  $\mathbf{X}$  son los vectores propios de  $\mathbf{C}_X = \frac{1}{n}\mathbf{X}\mathbf{X}^T$ .
- El valor  $i$ -ésimo de la diagonal de  $\mathbf{C}_Y$  es la varianza de  $\mathbf{X}$  a lo largo de  $\mathbf{p}_i$



Alaa Tharwat (2016). Principal component analysis - a tutorial

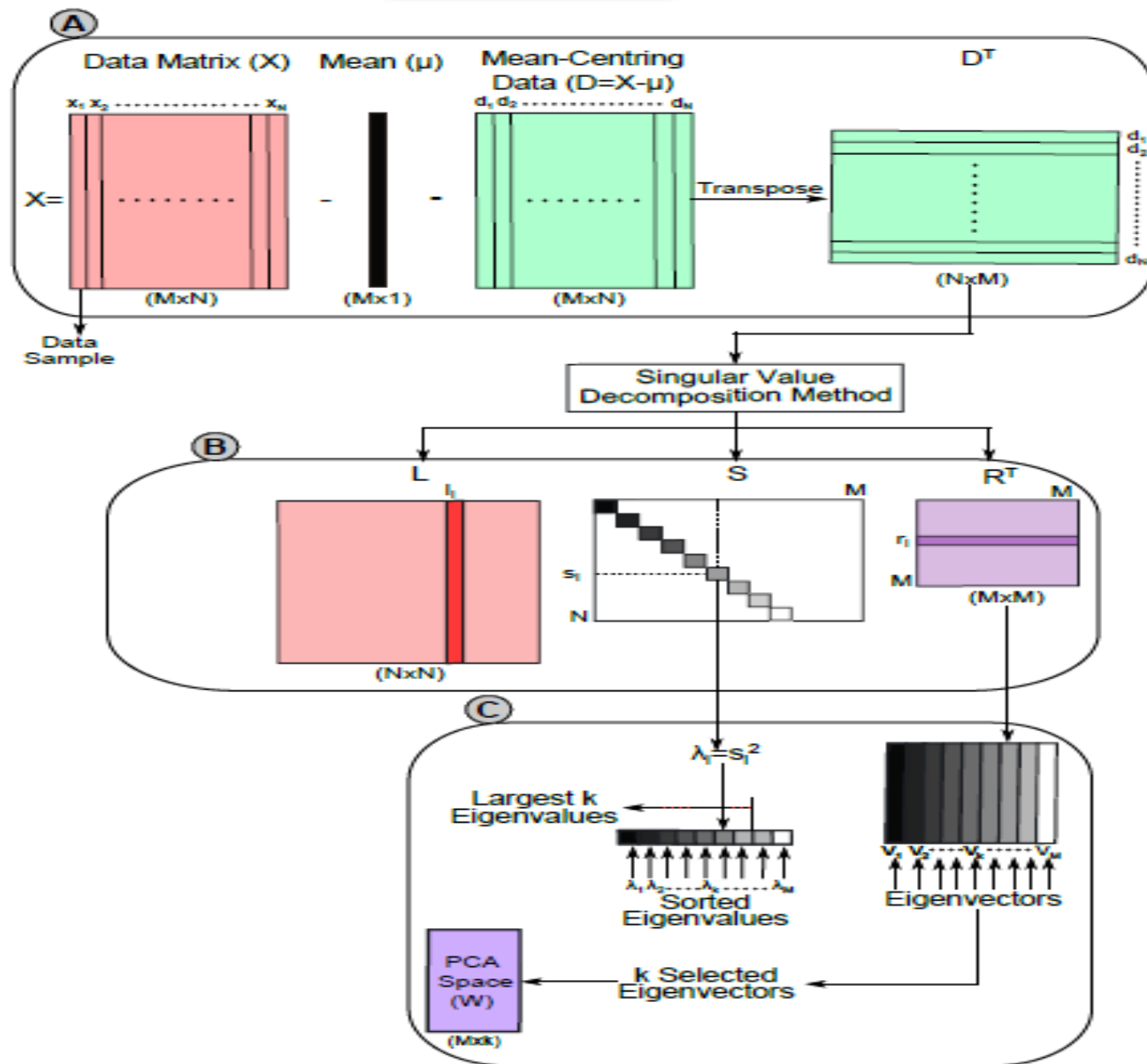
# SOLUCIÓN UTILIZANDO DESCOMPOSICIÓN EN VALORES SINGULARES (SVD)

OBJETIVO DE SVD:

Descomponer  $X$  en tres matrices

$$X = LSR^T = \begin{bmatrix} l_1 & \cdots & l_p \end{bmatrix} \begin{bmatrix} s_1 & 0 & 0 & 0 \\ 0 & s_2 & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & s_q \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} -r_1^T & - \\ -r_2^T & - \\ \vdots & \\ -r_q^T & - \end{bmatrix}$$

Alaa Tharwat (2016). Principal component analysis - a tutorial



Alaa Tharwat (2016).  
Principal component  
analysis - a tutorial