

# Reconocimiento Estadístico de Patrones - Tarea 1

Rubén Pérez Palacios Lic. Computación Matemática, Profesor: Johan Van Horebeek

February 9, 2023

## 1 Ejercicios

### 1.1 ChatGPT

#### 1.1.1 Problema



Para un estudio se mide la temperatura en diferentes partes del cuerpo de una muestra de personas. Un investigador expresa todas las temperaturas en grados Celsius.

Otro investigador convierte primero todas estas temperaturas a grados Fahrenheit. Si ambos deciden proyectar sus datos en la dirección de máxima varianza, ¿obtendrán los mismos resultados?



La proyección de los datos en la dirección de máxima varianza se realiza utilizando la técnica de análisis de componentes principales (PCA).

En el caso del estudio de las temperaturas en diferentes partes del cuerpo, ambos investigadores tienen la misma información, pero expresada en unidades diferentes (grados Celsius vs grados Fahrenheit).

La PCA se basa en la matriz de covarianza de los datos. Sin embargo, la matriz de covarianza solo mide las relaciones lineales entre variables, y no depende de las unidades en las que se expresan las variables. Por lo tanto, al calcular la matriz de covarianza de los datos en grados Celsius y grados Fahrenheit se obtendrán las mismas direcciones principales, independientemente de las unidades utilizadas.

Por lo tanto, si ambos investigadores proyectan sus datos en la dirección de máxima varianza utilizando PCA, obtendrán los mismos resultados, independientemente de si los datos están expresados en grados Celsius o grados Fahrenheit.

¿Cómo calificarías la respuesta de ChatGPT? ¿Cuál es tu respuesta? Demuéstrala formalmente.

#### 1.1.2 Respuestas

Aunque la respuesta y la intuición es correcta, no lo es la justificación.

A pesar de que la matriz de covarianza solo mide las relaciones lineales entre variables no necesariamente es independiente de como esten expresadas las variables. Puesto que estas pueden no ser una transformación lineal una de la otra, como es nuestro caso.

Sin embargo si se cumple en nuestro caso, veamos lo siguiente. Sean  $X$  y  $Y$  v.a. que miden la temperatura en difernetes partes del cuerpo en grados Celcius y grados Fahrenheit respectivamente. Recordemos que

$$Y = \frac{9}{5}X + 32.$$

Puesto que la varianza es invariante bajo translaciones y  $Var(aX) = a^2Var(X)$ .

$$Var(Y) = Var(\frac{9}{5}X + 32) = \frac{81}{25}Var(X),$$

es decir la varianza de  $Y$  es proporcional a la varianza de  $X$ , por lo tanto maximizar la varianza de  $X$  es análogo a maximizar la varianza de  $Y$ ,

$$\max_l \frac{Var(l^t X)}{||l||^2} \Longleftrightarrow \max_l \frac{Var(l^t Y)}{||l||^2},$$

con lo que concluimos que ambos obtendran la misma dirección de máxima varianza.

## 1.2 Normal Estandar

### 1.2.1 Problema

Supongamos que  $X = (X_1, X_2)$ ,  $Var(X_1) = Var(X_2) = 1$ ,  $E[X_1] = E[X_2] = 1$ ,  $yX_1 \perp X_2$ . Demuestra que cualquier dirección  $l$  da máxima varianza en las proyecciones.

### 1.2.2 Demostración

Por definición de  $Cov(X)$  obtenemos

$$Cov(X) = \begin{pmatrix} Var(X_1) & Cov(X_1, X_2) \\ Cov(X_2, X_1) & Var(X_2) \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix},$$

por lo tanto

$$\max_l \frac{l^t Cov(X) l}{||l||^2} = \max_l \frac{l^t l}{||l||^2} = \max_l 1,$$

lo cual es independiente de  $l$ . Con lo que concluimos cualquier dirección  $l$  da máxima varianza en las proyecciones.

### 1.3 Demuestra que $K = -\frac{1}{2}CD^2C$

Recordemos que

- $K$  : Es la matriz de producto punto entre variables, es decir  $K_{i,j} = X_i \dot{X}_j$ .
- $D^2$  : Es la matriz de distancias al cuadrado entre variables, es decir  $D_{i,j}^2 = d(X_i, X_j)^2$ .
- $C$  : Es la matriz identidad menos  $1/n$  sobre todas las entradas, es decir  $C = I - \frac{1}{n}J$ .

También en clase se demostro que

$$K = -\frac{1}{2}CD^2C$$

$$K_{i,j} = -\frac{1}{2} \left( D_{i,j}^2 - \frac{\sum_j D_{i,j}^2}{n} + \frac{\sum_j K_{i,j}}{n} - \frac{\sum_i D_{i,j}^2}{n} + \frac{\sum_i K_{i,j}}{n} \right)$$

$$\sum_i \sum_j D_{i,j}^2 = 2n \sum_i K_{i,i}$$

por lo que nos falta demostrar que

$$CD^2C = D_{i,j}^2 - \frac{1}{n} \sum_j D_{i,j}^2 - \frac{1}{n} \sum_i D_{i,j}^2 + \frac{1}{n^2} \sum_i \sum_j D_{i,j}^2.$$

Para ello usaremos la notación  $X_{i,\cdot}$  y  $X_{\cdot,j}$ , para denotar al vector renglón y columna respectivamente. Ahora veamos que

$$\begin{aligned}
(CD^2C)_{i,j} &= (CD^2)_{i,\cdot} C_{\cdot,j} && \text{Por definición de multiplicación de matrices} \\
&= \sum_k (CD^2)_{i,k} (C)_{k,j} && \text{Por definición de multiplicación de matrices} \\
&= \sum_k (C_{i,\cdot} D^2_{\cdot,k}) C_{k,j} && \text{Por definición de multiplicación de matrices} \\
&= \sum_k \left( \sum_h C_{i,h} D^2_{h,k} \right) C_{k,j} && \text{Por definición de multiplicación de matrices} \\
&= \sum_k \left( D^2_{i,k} - \sum_h \frac{1}{n} D^2_{h,k} \right) C_{k,j} && \text{Por definición de C} \\
&= \sum_k \left( D^2_{i,k} - \frac{1}{n} \sum_h D^2_{h,k} \right) C_{k,j} && \\
&= \left( D^2_{i,j} - \frac{1}{n} \sum_h D^2_{h,j} \right) + \sum_k -\frac{1}{n} \left( D^2_{i,k} - \frac{1}{n} \sum_h D^2_{h,k} \right) && \text{Por definición de C} \\
&= D^2_{i,j} - \frac{1}{n} \sum_h D^2_{h,j} - \frac{1}{n} \sum_k \left( D^2_{i,k} - \frac{1}{n} \sum_h D^2_{h,k} \right) && \\
&= D^2_{i,j} - \frac{1}{n} \sum_i D^2_{i,j} - \frac{1}{n} \sum_j D^2_{i,j} + \frac{1}{n^2} \sum_i \sum_j D^2_{i,j} && \\
& && (10)
\end{aligned}$$

Por lo tanto concluimos la demostración.

## 1.4 Plano de direcciones PCA

### 1.4.1 Problema

Revisa el video sobre la maximización del cociente de Rayleigh: <https://youtu.be/8TBpSUXcDww>.

Haz unos pequeños cambios necesarios para demostrar que el segundo vector propio de  $\text{Cov}(X)$  es la solución del problema de maximizar el cociente bajo la restricción adicional de ser ortogonal al primer vector propio.

### 1.4.2 Demostración

Primero recordemos lo siguiente de la demostración. Sea  $A$  una matriz simétrica entonces

- Su SVD es  $A = U\Lambda U^t$ , donde  $U^t = U^{-1}$  y  $\Lambda$  es una matriz diagonal con los valores propios de  $A$ .
- Definimos a  $\Lambda^{1/2} := \text{Diag}(\sqrt{\mu_i})$ , donde  $\mu_i$  son los valores propios de  $A$ , entonces podemos ver que  $\Lambda = \Lambda^{1/2}\Lambda^{1/2}$  y por lo tanto  $(U\Lambda^{1/2}U^t)(U\Lambda^{1/2}U^t) = A^{1/2}A^{1/2} = A$ .

- El vector dirección que maximiza la varianza es  $U^t e_1 = v_1$  el primer vector propio de  $Cov(X)$  con valor  $\mu_1$ .
- Debido a que  $U$  es una matriz ortonormal entonces conserva ortogonalidades entonces  $l_1^t l = 0$  ssi  $(U^t l_1)^t (U^t l) = 0$  ssi  $e_1^t y = 0$ .
- Por último, si  $v^t e_1 = 0$  es ssi  $\sum v_i e_i = 0$  lo cual es ssi  $v_1 = 0$ .

Entonces podemos ver que

$$\max_{l: l_1^t l = 0} \frac{l^t Cov(X) l}{l^t l} = \max_{l: l_1^t l = 1} \frac{l^t (U \Lambda^{1/2} U^t) (U \Lambda^{1/2} U^t) l}{l^t l} \quad \text{Puesto que } Cov(X) \text{ es simétrica} \quad (11)$$

$$= \max_{l: l_1^t l = 0} \frac{l^t (U \Lambda^{1/2} U^t) (U \Lambda^{1/2} U^t) l}{l^t (U U^t) l} \quad \text{Ya que } U \text{ es una matriz unitaria} \quad (12)$$

$$= \max_{y: y^t e_1 = 0} \frac{y \Lambda y^t}{y^t y} \quad \text{Con el cambio de variable } y = U^t l \quad (13)$$

$$= \max_{y: y^t e_1 = 0} \frac{\sum_{i=1}^n \mu_i y_i^2}{\sum_{i=1}^n y_i^2} \quad \text{Por definición de } \Lambda \quad (14)$$

$$= \max_y \frac{\sum_{i=2}^n \mu_i y_i^2}{\sum_{i=2}^n y_i^2} \quad \text{Ya que } y^t e_1 = 0 \text{ entonces } y_1 = 0 \quad (15)$$

$$\leq \max_y \frac{\sum_{i=1}^n \mu_2 y_i^2}{\sum_{i=1}^n y_i^2} \quad \text{Debido a que } \mu_2 \geq \mu_i, \forall i > 2 \quad (16)$$

$$= \max_y \mu_2 \quad (17)$$

$$(18)$$

Ahora veamos que si  $y = (0, 1, \dots, 0)$  entonces

$$\frac{\sum \mu_i y_i^2}{\sum y_i^2} = \frac{\mu_2}{1} = \mu_2.$$

Por lo tanto concluimos que  $v_2$  el segundo vector propio de  $Cov(X)$ , es la dirección que maximiza la varianza y que además es ortogonal a  $v_1$ , cuyo máximo es  $\mu_2$ .

## 1.5 Visualización de Datos

### 1.5.1 Problema

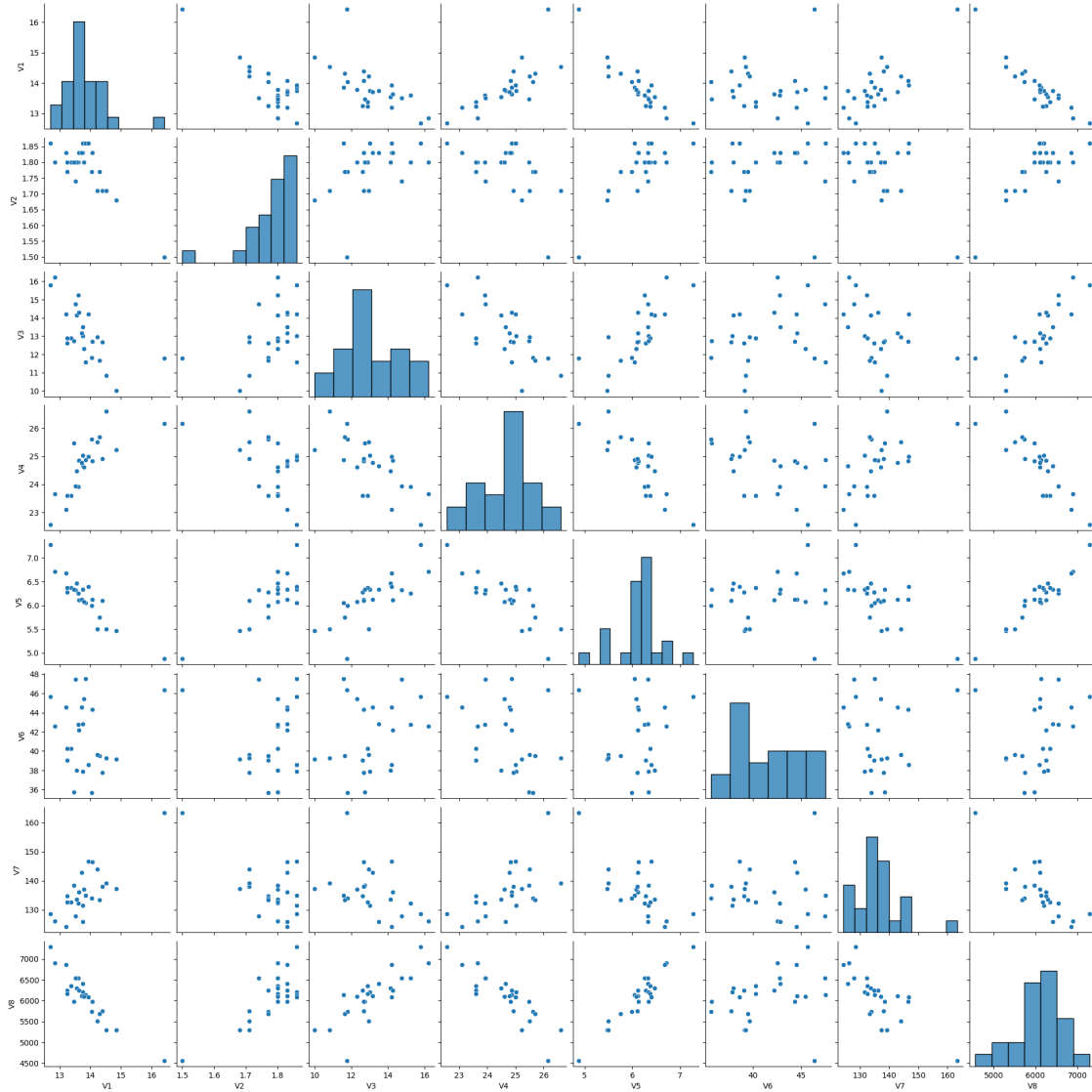
En el archivo heptatlon se pueden consultar los tiempos y el puntaje final (score) de 25 atletas que participaron en el heptatlon durante los juegos olímpicos de Seoul.

- Busca visualizaciones informativas de estos datos multivariados.
- Haz un análisis de componentes principales con los tiempos (sinscore). Hay una relación entre el score y las proyecciones sobre el primer CP? Puedes leer los datos con: `dj-read.table("heptatlon")`

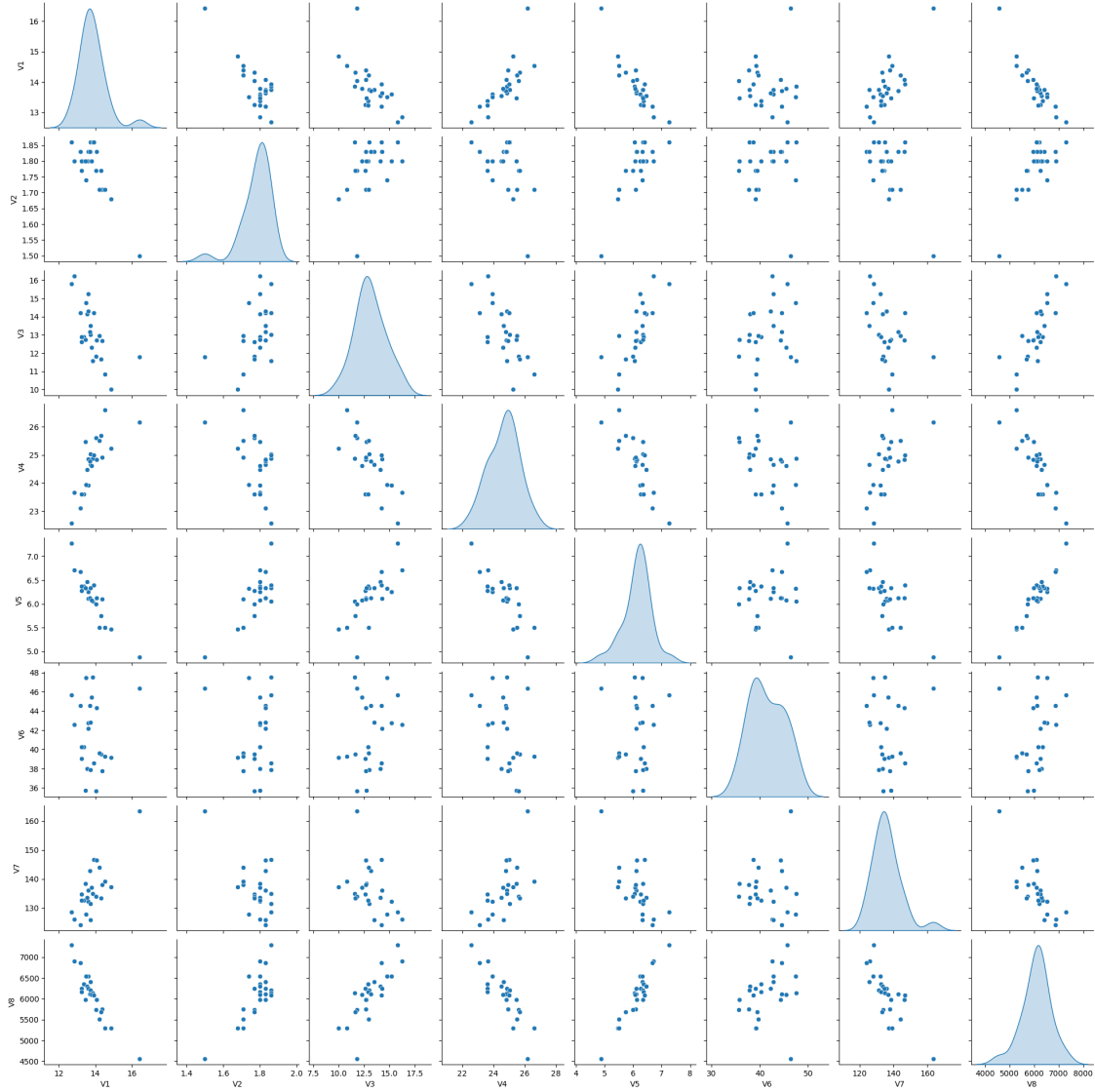
### 1.5.2 Solución

- a) Renombramos (hurdles,highjump,shot,run200m,longjump,javelin,run800m,score) como (V1, V2, V3, V4, V5, V6, V7, V8)

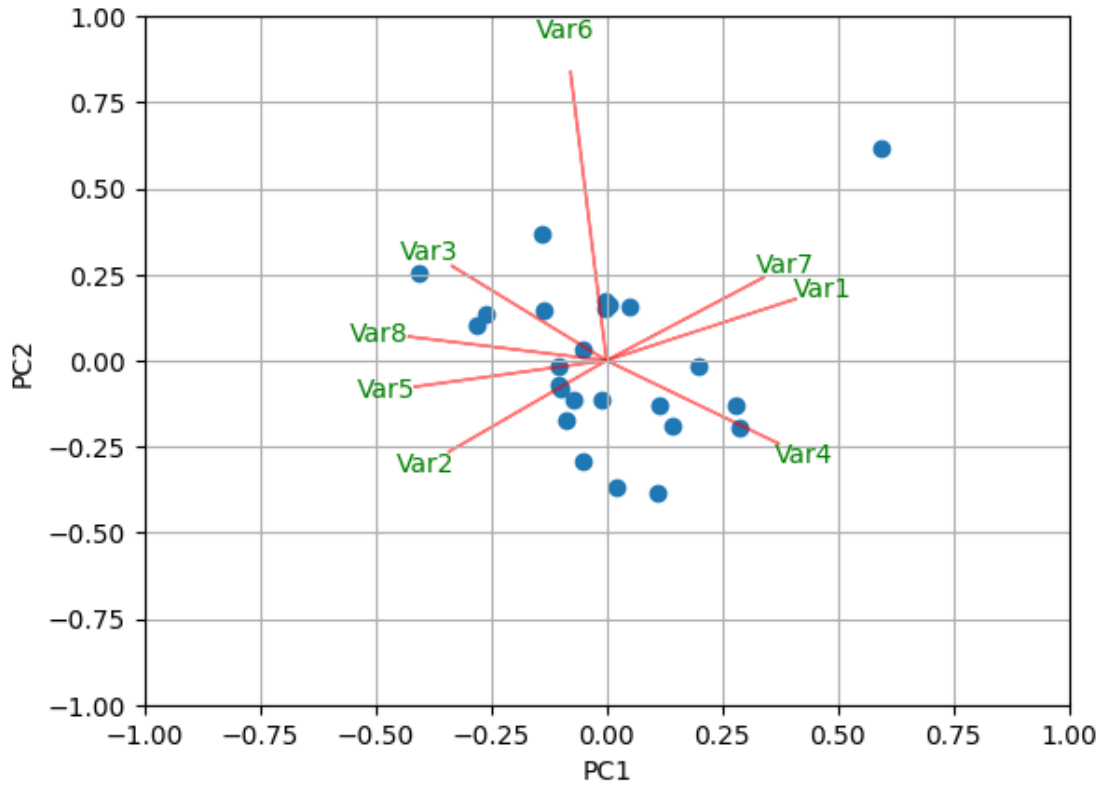
Distribuciones conjuntas a pares de  $(V_i, V_j)$  con  $i \neq j$  y distribución marginal discreta  $V_i$  cuando  $i = j$ .



Distribuciones conjuntas a pares de  $(V_i, V_j)$  con  $i \neq j$  y distribución marginal continua  $V_i$  cuando  $i = j$ .



Podemos ver en el biplot que  $v_7$ ,  $v_1$  y  $v_4$  están fuertemente relacionada con el PC1 e inversamente relacionadas  $v_2$ ,  $v_5$ ,  $v_8$ ,  $v_3$ . En cambio la  $v_6$  está fuertemente relacionada con PC2. Así mismo podemos ver que la hay bastante correlación entre  $v_7$ ,  $v_1$  y  $v_4$ , así como bastante correlación entre  $v_2$ ,  $v_5$ ,  $v_8$ ,  $v_3$ .



b)

[ 1.19624885e-03, -1.05204013e-04, -2.09849444e-03, 1.47499229e-03, -7.92784821e-04, -1.57865298e-03, 1.12710147e-02, -9.99930909e-01]

[ 1.92835321e-02, -6.05680149e-04, 8.27323011e-02, -1.73321242e-02, 1.08856631e-03, 3.27484685e-01, 9.40817259e-01, 9.91075581e-03]

[-4.42281031e-02, 5.20709703e-03, 5.23034341e-02, 2.52425088e-02, 3.33247513e-02, -9.42375189e-01, 3.24711313e-01, 4.99545357e-03]

[ 1.27735048e-01, -3.39503633e-02, 9.78341573e-01, 1.16494549e-01, -5.42199304e-02, 1.74323988e-02, -9.25021755e-02, -2.75215967e-03]

[ 2.99253782e-02, -6.54426901e-02, 1.10968349e-01, -9.90594181e-01, 8.48590375e-03, -2.83264257e-02, -1.87933943e-02, -1.82525884e-03]

[-9.77750566e-01, 1.11049607e-02, 1.19861716e-01, -1.94015933e-02, -1.64765138e-01, 4.41739357e-02, -6.02027169e-03, -1.45801494e-03]

[ 1.56298683e-01, 7.45409331e-02, -6.80079688e-02, -1.54259545e-02, -9.81458446e-01, -3.97783256e-02, 1.75083617e-02, 1.33740223e-03]