

Tarea 01 - Reconocimiento de patrones
Giovanni Gamaliel López Padilla

Problema 2

Supongamos que $X=(X_1, X_2)$, $\text{Var}(X_1)=\text{Var}(X_2)=1$.

a) Supongamos que X_1 y X_2 son v.a. independientes con promedio 0. Verifica que cualquier dirección l da máxima varianza en las proyecciones.

Como X_1 y X_2 son v.a independientes entonces, la matriz de covarianza $\text{Cov}(X)$ es:

$$\text{Cov}(X) = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

por lo que se obtiene que $\text{Cov}(X) = \mathbb{I}$. Entonces, se obtiene que:

$$\begin{aligned} \max_{\|l\|} \frac{l^t \text{Cov}(X) l}{l^t l} &= \max_{\|l\|} \frac{l^t \mathbb{I} l}{l^t l} \\ &= \max_{\|l\|} \frac{l^t l}{l^t l} \\ &= \max_{\|l\|} 1 \\ &= 1 \end{aligned}$$

por lo tanto, se maximiza la varianza para cualquier dirección de l en las proyecciones.

b) Supongamos que X_1 y X_2 son v.a. dependientes. Calcula la primer componente principal a mano. ¿Qué particularidad tiene?

Suponiendo de la covarianza de X_1 Y X_2 es a , entonces, la matriz de covarianza es:

$$\text{Cov}(X) = \begin{pmatrix} 1 & a \\ a & 1 \end{pmatrix}$$

Calculando la primer componente l , se obtiene que los valores propios de $\text{Cov}(X)$ es:

$$\begin{aligned} |\text{Cov}(X) - \lambda \mathbb{I}| &= 0 \\ \begin{vmatrix} 1 - \lambda & a \\ a & 1 - \lambda \end{vmatrix} &= 0 \\ (1 - \lambda)^2 - a^2 &= 0 \\ (1 - \lambda - a)(1 - \lambda + a) &= 0 \\ \lambda_1 &= 1 - a \\ \lambda_2 &= 1 + a \end{aligned}$$

Suponiendo que $a > 0$, entonces λ_2 es el eigenvalor mayor. Calculando los vectores propios relacionados a λ_2 , se obtiene que:

$$\begin{aligned} \begin{pmatrix} -a & a \\ a & -a \end{pmatrix} &= \begin{pmatrix} 0 \\ 0 \end{pmatrix} \\ \begin{pmatrix} -1 & 1 \\ a & -a \end{pmatrix} &= \begin{pmatrix} 0 \\ 0 \end{pmatrix} \\ \begin{pmatrix} -1 & 1 \\ 0 & 0 \end{pmatrix} &= \begin{pmatrix} 0 \\ 0 \end{pmatrix} \\ -x_1 + x_2 &= 0 \\ x_1 &= x_2 \end{aligned}$$

por lo tanto, el vector propio asociado a λ_2 es $v_2 = [x_1, x_1]^T$. La particularidad que tiene es que las componentes no tienen un valor determinado por lo que es necesario elegir el parámetro x_1 y en seguida normalizar el vector.

Problema 3

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 primer vector propio.

Al final del video se obtiene que una solución del problema descrito en la ecuación 1 es el primer vector propio de $\text{Cov}(X)$.

$$\max_{||l||} \frac{l^t \text{Cov}(X) l}{l^t l} \quad (1)$$

Realizando un cambio de base a la ecuación 1, se obtiene la ecuación 2.

$$\max_{||y||} \frac{y^t \Lambda y}{y^t y} \rightarrow \max_{||y||} \frac{\sum_i \mu_i y_i^2}{\sum_i y_i^2} \quad (2)$$

Tomando en cuenta que $\mu_1 \geq \mu_2 \geq \dots \geq \mu_i$, donde $y_1 = (1, 0, \dots, 0)$. Entonces se propone que $y_2 = (0, 1, 0, \dots, 0)$, esto con el propósito que y_1 y y_2 sean ortogonales. Usando y_2 en la ecuación 2, se obtiene que una solución es μ_2 , el cual es el segundo eigenvalor de la matriz $\text{Cov}(X)$. Devolviendo a la base original a y_2 se obtiene que:

$$\begin{aligned} y_2 &= U^t l_2 \\ U y_2 &= U U^t l_2 \\ U y_2 &= l_2 \\ u_2 &= l_2 \end{aligned}$$

donde u_2 es el segundo eigenvector de la matriz $\text{Cov}(X)$ el cual es ortogonal a u_1 .

Problema 4

Considera los datos [oef2.data](#). Se trata de los promedios mensuales de la temperatura (en Celsius) en 35 estaciones canadienses de monitoreo. El interés es comparar las estaciones entre sí en base de sus curvas de temperatura. Considerando las 12 mediciones por estación como un vector X , aplica un análisis de componentes principales. Interpreta y dibuja (como curva) los primeros dos componentes, p_1 , p_2 , es decir grafica $\{(i, p_{1i})\}$ y $\{(i, p_{2i})\}$. Agrupa e interpreta las estaciones en el biplot (ten en mente un mapa de Canada).

En la figura 1 se muestra las gráficas lineales de los valores de cada componente. Se observa que el comportamiento es semejante en las dos componentes. La primer componente conserva una menor variación en sus valores a comparación de la segunda componente.

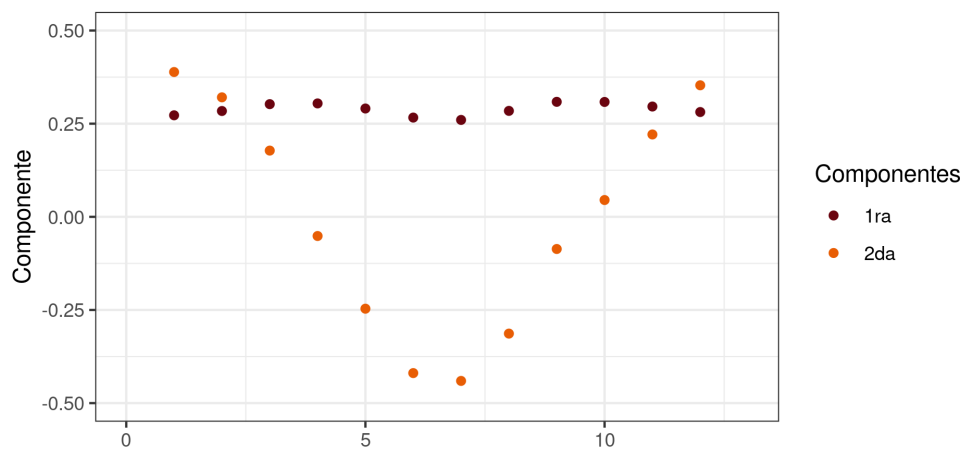
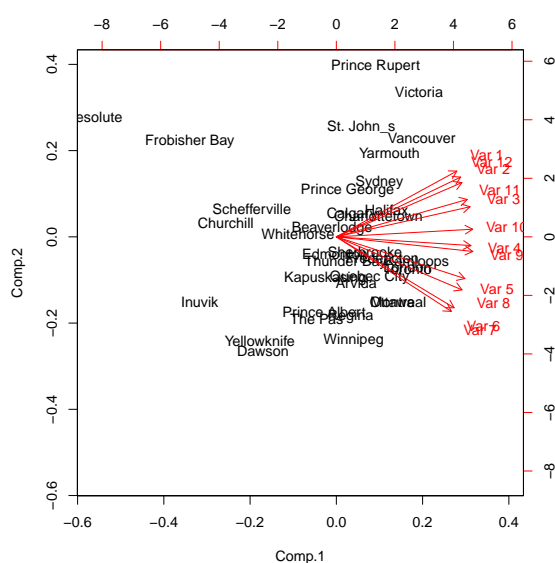


Figura 1: Valores de la primer y segunda componente obtenida con PCA.



(a) Biplot de la primer y segunda componente obtenido de los datos del archivo [oef2.data](#).



(b) Mapa de Canada. Obtenido de [Geology.com](#)

Figura 2

En la figura 2a se muestra la gráfica de biplots obtenida al aplicar PCA a los datos del archivo [oef2.data](#). Se observa como la distribución de valores se concentra para las estaciones que tienen una latitud semejante (figura 2b). En cambio para la estación ubicada en Resolute se encuentra lejana de las demás. Esto puede dar la interpretación que mientras más se encuentren en la izquierda de la gráfica las estaciones reportaran una temperatura menor. En caso contrario, estas reportaran una temperatura mayor.

Problema 5

En este ejercicio usamos resultados del heptatlón femenina de los pasados juegos olímpicos de Tokyo (2021). En el archivo [heptatlonTokyo](#) se pueden consultar los tiempos/distancias y el puntaje final (score) de 20 atletas.

- a) Describe de manera general los datos sin considerar la columna con los puntajes finales, usando visualizaciones ilustrativas. Toma en cuenta que son pocas observaciones. Así, no será posible llegar a conclusiones fuertes.
- b) Un problema en un heptatlón es cómo convertir los resultados obtenidos en las diferentes pruebas en un puntaje final. Explora la utilidad de PCA usando la proyección de los resultados de 4 las pruebas sobre el primer componente como una alternativa al puntaje final. ¿Cómo se relaciona con el puntaje final oficial? Información sobre cómo se calcula actualmente el puntaje: [Puntaje heptatlon](#).