

Capítulo 2

Análisis en Componentes Principales (ACP).

Como antes se ha dicho, el Análisis en Componentes Principales, ACP, considera una matriz R de datos iniciales de carácter no simétrico: Sus componentes, r_{ij} , son valores numéricos que, de manera continua, miden la variable j -ésima (columnas) en el i -ésimo individuo (filas).

Veamos a continuación cómo el Análisis Factorial General se adapta a esta situación, de manera que tengamos en cuenta que se trata de una tabla estadística y no de una mera matriz numérica R que hubiera que reducir o explicar en base a un conjunto menor de números.

Queda pues concretada la situación inicial en estos términos:

1. El espacio de las filas está constituido por n individuos.
2. El espacio de las columnas está constituido por p variables.

Así pues, disponemos de una tabla de datos $R = (r_{ij})$; $i = 1, \dots, n$; $j = 1, \dots, p$; sobre la que vamos a efectuar en primer lugar el Análisis en \mathbb{R}^p , previa concreción de la matriz X .

2.1. Análisis en el espacio de las variables \mathbb{R}^p .

Según el AFG, consideraremos en una primera fase los n puntos-fila en el espacio vectorial \mathbb{R}^p de las columnas (variables). En primer lugar, y dada la naturaleza estadística de la tabla inicial R , vamos a definir sobre qué matriz, X , transformada de la tabla R , vamos a realizar el AFG.

En realidad, cuando se tiene una situación práctica del tipo que estamos considerando (variables-individuos), las variables pueden ser muy heterogéneas en cuanto a sus valores medios. También desde luego lo pueden ser respecto a su dispersión, provocada por escalas de medida a veces muy diferentes, pero esta otra fuente de heterogeneidad no la consideramos en una primera fase. Así pues, ¿cómo eliminar el efecto de la heterogeneidad de medidas?

Hagamos una interpretación geométrica de la situación. No olvidemos que hemos de realizar (Análisis en \mathbb{R}^p) el ajuste en \mathbb{R}^p (espacio de columnas-variables) de un subespacio que describa aproximadamente la nube de puntos-fila (los individuos), con los criterios definidos en el AFG.

Este subespacio, como tal, contiene siempre al origen en el planteamiento básico del AFG, en donde no se considera la naturaleza estadística de la tabla sino meramente su carácter numérico. Ahora, sin embargo, tenemos una información adicional sobre los valores numéricos, información estadística, que permite tratar la heterogeneidad de las medidas de las variables como un factor tanto a considerar como a eliminar, en su posible efecto no deseado, a la hora de describir la proximidad de los puntos fila.

Gráficamente lo que ocurre es que la posición de la nube de puntos-fila (individuos) cambia, respecto del origen, frente a la adoptada en el AFG en donde no se considera el efecto de las medias en cada variable.

En lenguaje geométrico, es claro que el espacio afín \mathbb{H}_1 es presumiblemente más apropiado para describir la nube de puntos-fila que el espacio vectorial \mathbb{H}_0 (subespacio del \mathbb{R}^p). Considerar este subespacio afín, no

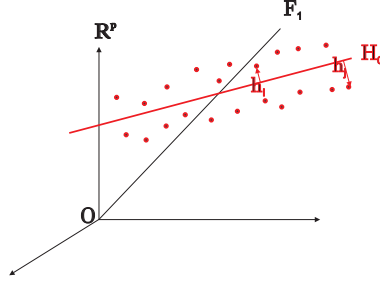


Figura 2.1: Recta afín

altera la filosofía general del AFG, ya que lo que estadísticamente interesa es la forma de la nube, no la posición de ésta respecto del origen.

Técnicamente hablando, interesa pasar el origen al centro de gravedad (punto medio, cuyas coordenadas son las dadas por las p medias, \bar{r}_j , en las p variables), porque así, la búsqueda de la recta \mathbb{H}_1 (subespacio afín, cuando el origen está en 0) puede reducirse a la búsqueda de un \mathbb{H}_0 tal como se ha realizado en el AFG, pero respecto de un origen colocado en el centro de gravedad.

Analíticamente, estas consideraciones geométricas equivalen a lo siguiente: Sean h_i y h_j las proyecciones de dos puntos sobre una recta \mathbb{H}_1 . Consideremos el criterio de maximizar la suma $\sum_{i,j} (h_i - h_j)^2$, es decir maximizar la suma de todas las diferencias de proyecciones de todos los pares al cuadrado. Si esta suma se desarrolla, se puede ver que vale:

$$2n \sum_{i=1}^n (h_i - \bar{h})^2 \quad (2.1)$$

en donde \bar{h} denota la proyección media, es decir la media aritmética de las proyecciones de todos los puntos sobre la recta \mathbb{H}_1 . En efecto

$$\begin{aligned} \sum_{i,j=1}^n (h_i - h_j)^2 &= \sum_{i,j=1}^n [h_i - \bar{h} - \bar{h}_j + \bar{h}]^2 = \\ &= \sum_{i,j=1}^n [(h_i - \bar{h}) - (h_j - \bar{h})]^2 = \\ &= \sum_{i,j=1}^n (h_i - \bar{h})^2 + \sum_{i,j=1}^n (h_j - \bar{h})^2 - 2 \sum_{i,j=1}^n (h_i - \bar{h})(h_j - \bar{h}) = \\ &= 2n \sum_{i=1}^n (h_i - \bar{h})^2 - 2 \sum_{i,j=1}^n (h_i - \bar{h})(h_j - \bar{h}) = 2n \sum_{i=1}^n (h_i - \bar{h})^2 \end{aligned}$$

El último paso, se ha dado en virtud de que:

$$\begin{aligned} \sum_{i,j=1}^n (h_i - \bar{h})(h_j - \bar{h}) &= \sum_{i,j=1}^n h_i h_j - \bar{h} \sum_{i,j=1}^n h_i - \bar{h} \sum_{i,j=1}^n h_j + \sum_{i,j=1}^n \bar{h}^2 = \\ &= \sum_{i,j=1}^n h_i h_j - 2\bar{h} \sum_{i,j=1}^n h_i + n^2 \bar{h}^2 = \\ &= n \sum_{i=1}^n h_i \left(\frac{\sum_{j=1}^n h_j}{n} \right) - 2n\bar{h} \sum_{i=1}^n \frac{h_i}{n} + n^2 \bar{h}^2 = \end{aligned}$$

$$\begin{aligned}
&= n^2 \bar{h} \left(\frac{\sum_{i=1}^n h_i}{n} \right) - 2n^2 \bar{h}^2 + n^2 \bar{h}^2 = \\
&= n^2 \bar{h}^2 - 2n^2 \bar{h}^2 + n^2 \bar{h}^2 = 0
\end{aligned}$$

C.Q.D.

Es claro entonces, que si hicieramos $\bar{h} = 0$, el problema de maximización planteado, que conducirá a la recta \mathbb{H}_1 , sería un caso particular del problema de maximización planteado en el AFG (maximizar la suma de cuadrados de todas las proyecciones). Pero ¿qué es hacer $\bar{h} = 0$? La respuesta es evidente: *Es trasladar el origen de coordenadas al centro de gravedad de la nube.* Es decir, considerar los valores

$$x_{ij}^* = r_{ij} - \bar{r}_j$$

La matriz final X , no obstante, se define así:

$$x_{ij} = \frac{r_{ij} - \bar{r}_j}{\sqrt{n}} \quad (2.2)$$

transformación pues que pasa de la matriz inicial $R = (r_{ij})$ a la $X = (x_{ij})$. El factor $\frac{1}{\sqrt{n}}$ se añade para que, al aplicar el AFG, y tener que considerar $X'X$, aparezca la matriz de covarianzas muestrales. En efecto,

$$X'X = \left(\frac{r_{ij} - \bar{r}_j}{\sqrt{n}} \right)'_{p \times n} \left(\frac{r_{ij} - \bar{r}_j}{\sqrt{n}} \right)_{n \times p} = C_{p \times p} \quad (2.3)$$

es la matriz de covarianzas muestrales. Es decir,

$$X = (x_{ij})_{n \times p} = \begin{pmatrix} \frac{r_{11} - \bar{r}_1}{\sqrt{n}} & \frac{r_{12} - \bar{r}_2}{\sqrt{n}} & \dots & \frac{r_{1p} - \bar{r}_p}{\sqrt{n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{r_{n1} - \bar{r}_1}{\sqrt{n}} & \frac{r_{n2} - \bar{r}_2}{\sqrt{n}} & \dots & \frac{r_{np} - \bar{r}_p}{\sqrt{n}} \end{pmatrix}_{n \times p}$$

de modo que:

$$\begin{aligned}
X'_{p \times n} X_{n \times p} &= \begin{pmatrix} \frac{r_{11} - \bar{r}_1}{\sqrt{n}} & \dots & \frac{r_{n1} - \bar{r}_1}{\sqrt{n}} \\ \vdots & \ddots & \vdots \\ \frac{r_{1p} - \bar{r}_p}{\sqrt{n}} & \dots & \frac{r_{np} - \bar{r}_p}{\sqrt{n}} \end{pmatrix}_{p \times n} \begin{pmatrix} \frac{r_{11} - \bar{r}_1}{\sqrt{n}} & \dots & \frac{r_{1p} - \bar{r}_p}{\sqrt{n}} \\ \vdots & \ddots & \vdots \\ \frac{r_{n1} - \bar{r}_1}{\sqrt{n}} & \dots & \frac{r_{np} - \bar{r}_p}{\sqrt{n}} \end{pmatrix}_{n \times p} = \\
&= \begin{pmatrix} \sum_{i=1}^n \frac{(r_{i1} - \bar{r}_1)^2}{n} & \sum_{i=1}^n \frac{(r_{i1} - \bar{r}_1)(r_{i2} - \bar{r}_2)}{n} & \dots \\ \sum_{i=1}^n \frac{(r_{i2} - \bar{r}_2)(r_{i1} - \bar{r}_1)}{n} & \dots & \dots \\ \vdots & \ddots & \vdots \\ \dots & \dots & \dots \end{pmatrix}_{p \times p}
\end{aligned}$$

En una segunda fase, como antes decíamos, también debe tenerse en cuenta la *heterogeneidad* de las escalas de medida, que produce en muchos casos una dispersión muy heterogénea en las variables. La forma típica de reducir el efecto de la heterogeneidad en las dispersiones de las variables, a la hora de medir la proximidad de los individuos, es efectuar el paso a una matriz X dada en la forma:

$$X = (x_{ij}) = \left(\frac{r_{ij} - \bar{r}_j}{s_j \sqrt{n}} \right); s_j^2 = \frac{1}{n} \sum_{i=1}^n (r_{ij} - \bar{r}_j)^2 \quad (2.4)$$

Dados dos individuos (i -ésimo e i' -ésimo), la distancia al cuadrado entre ellos será (en sentido euclídeo clásico):

$$\begin{aligned}
d^2(i, i') &= \sum_{j=1}^p \left(\frac{r_{ij} - \bar{r}_j}{s_j \sqrt{n}} - \frac{r_{i'j} - \bar{r}_j}{s_j \sqrt{n}} \right)^2 = \\
&= \sum_{j=1}^p \frac{(r_{ij} - r_{i'j})^2}{s_j^2 \cdot n} = \frac{1}{n} \sum_{j=1}^p \frac{(r_{ij} - r_{i'j})^2}{s_j^2}
\end{aligned}$$

y es evidente que, al dividir cada término por la varianza muestral s_j^2 correspondiente de cada variable ($j = 1, \dots, p$), se reduce el efecto de la heterogeneidad de las mismas, de modo que cada variable aporta una contribución análoga a la distancia o proximidad entre individuos. En este caso, además, es claro que la matriz $X'X$ no es otra cosa que la matriz de correlaciones muestrales. En efecto, obsérvese que, a partir de la matriz C dada en la expresión 2.3, la expresión 2.4 puede escribirse matricialmente en la forma $V^{-\frac{1}{2}}CV^{-\frac{1}{2}}$, en donde $V = \text{diag}(s_1^2, s_2^2, \dots, s_p^2)$. Pero $V^{-\frac{1}{2}}CV^{-\frac{1}{2}} = X'X = \mathbf{R}$ en donde $\mathbf{R}_{p \times p}$ es la matriz de coeficientes de correlación tipo Pearson.

2.1.1. Resumen.

El ajuste a la nube de puntos-fila (individuos) en \mathbb{R}^p (espacio de las variables), en el Análisis de Componentes Principales, se efectúa a dos niveles:

1. Sobre la matriz transformada: $X = \left(\frac{r_{ij} - \bar{r}_j}{\sqrt{n}} \right)_{n \times p}$ (*Análisis en Componentes Principales*), de modo que se aplica el AFG a $(X'_{p \times n} X_{n \times p})_{p \times p}$ para ajustar la nube en \mathbb{R}^p .
2. Sobre la matriz transformada: $X = \left(\frac{r_{ij} - \bar{r}_j}{s_j \sqrt{n}} \right)$ (*Análisis en Componentes Principales Normalizado*) de modo que se aplica el AFG a $(X'X)$ que es la matriz de correlaciones muestrales.

Las coordenadas de los puntos-fila sobre los ejes factoriales soportes de los autovectores u_α de $X'X$ son globalmente expresadas como Xu_α , y este vector columna $n \times 1$ tiene en cada componente la coordenada de un punto-fila respecto del α -ésimo eje. Por otra parte, según vimos en AFG, $Xu_\alpha = \sqrt{\lambda_\alpha} v_\alpha$, en donde v_α es el autovector α -ésimo en el otro espacio \mathbb{R}^n ; es decir, autovector de XX' .

2.2. Ajuste en \mathbb{R}^n de la nube de puntos-columna.

Como en los planteamientos generales del AFG, también en el ACP cabe ajustar a los p puntos-columna (variables), en \mathbb{R}^n , un subespacio vectorial o afín. Lo que ocurre en ACP es que la matriz inicial R , o sus transformadas X , no son simétricas en i, j . Es más, la transformación que lleva de R a X (en los dos casos antes definidos) no es simétrica en los índices i, j , y ha de interpretarse en un sentido muy diferente.

La transformación $x_{ij} = r_{ij} - \bar{r}_j$ implica trasladar el origen en \mathbb{R}^p , al punto *centro de gravedad* de la nube de los n puntos-fila. En cambio, interpretada en \mathbb{R}^n (es decir en j), la transformación $x_{ij} = r_{ij} - \bar{r}_j$ equivale a efectuar la operación matricial $X = P \cdot R$, en donde

$$P = (p_{ij}); \quad p_{ij} = \begin{cases} 1 - \frac{1}{n} & i = j \\ -\frac{1}{n} & i \neq j \end{cases}$$

Comprobemos que, en efecto, es $X = P \cdot R$.

$$P \cdot R = \begin{pmatrix} 1 - \frac{1}{n} & -\frac{1}{n} & \cdots & -\frac{1}{n} \\ -\frac{1}{n} & 1 - \frac{1}{n} & \cdots & -\frac{1}{n} \\ \vdots & \vdots & \ddots & \vdots \\ -\frac{1}{n} & -\frac{1}{n} & \cdots & 1 - \frac{1}{n} \end{pmatrix}_{n \times n} \begin{pmatrix} r_{11} & r_{12} & \cdots & r_{1p} \\ r_{21} & r_{22} & \cdots & r_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ r_{n1} & r_{n2} & \cdots & r_{np} \end{pmatrix}_{n \times p} =$$

$$\begin{aligned}
&= \begin{pmatrix} r_{11} - \frac{r_{11}}{n} - \frac{r_{21}}{n} - \dots - \frac{r_{n1}}{n} & r_{12} - \frac{r_{12}}{n} - \frac{r_{22}}{n} - \dots - \frac{r_{n2}}{n} & \dots \\ r_{21} - \frac{r_{21}}{n} - \frac{r_{22}}{n} - \dots - \frac{r_{2n}}{n} & \dots & \dots \\ \vdots & \ddots & \vdots \\ r_{n1} - \frac{r_{n1}}{n} - \frac{r_{n2}}{n} - \dots - \frac{r_{nn}}{n} & \dots & \dots \end{pmatrix} = \\
&= \begin{pmatrix} r_{11} - \bar{r}_1 & r_{12} - \bar{r}_2 & \dots & r_{1p} - \bar{r}_p \\ r_{21} - \bar{r}_1 & r_{22} - \bar{r}_2 & \dots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ r_{n1} - \bar{r}_1 & r_{n2} - \bar{r}_2 & \dots & r_{np} - \bar{r}_p \end{pmatrix} = X
\end{aligned}$$

Esta transformación —que también aparece, como es sabido, en el estudio de los Modelos Lineales— significa geoméricamente algo bien distinto de una traslación del origen en \mathbb{R}^p : Equivale a una proyección paralela a la primera bisectriz en \mathbb{R}^n .

En el caso del ACP normalizado, cuando se efectúa además el cambio de escala $\frac{r_{ij} - \bar{r}_j}{s_j \sqrt{n}}$ en \mathbb{R}^p , dividiendo cada coordenada de un punto-fila por $s_j \sqrt{n}$, se realiza una deformación de la nube de puntos-columna (p en total) de tal manera que se llevan los p puntos a una distancia unidad del origen. En efecto:

$$\begin{aligned}
d^2(j, 0) &= \sum_{i=1}^n \left(\frac{r_{ij} - \bar{r}_j}{s_j \sqrt{n}} - 0 \right)^2 = \frac{1}{n} \sum_{i=1}^n \frac{(r_{ij} - \bar{r}_j)^2}{s_j^2} = \\
&= \frac{1}{s_j^2} \left(\frac{1}{n} \sum_{i=1}^n (r_{ij} - \bar{r}_j)^2 \right) = 1.
\end{aligned}$$

Y la distancia entre dos puntos-columna, j y j' , será:

$$\begin{aligned}
d^2(j, j') &= \sum_{i=1}^n \left(\frac{r_{ij} - \bar{r}_j}{s_j \sqrt{n}} - \frac{r_{ij'} - \bar{r}_{j'}}{s_{j'} \sqrt{n}} \right)^2 = \\
&= \frac{1}{n} \sum_{i=1}^n \frac{(r_{ij} - \bar{r}_j)^2}{s_j^2} + \frac{1}{n} \sum_{i=1}^n \frac{(r_{ij'} - \bar{r}_{j'})^2}{s_{j'}^2} - \frac{2}{n} \sum_{i=1}^n \frac{(r_{ij} - \bar{r}_j)(r_{ij'} - \bar{r}_{j'})}{s_j s_{j'}} = \\
&= \frac{s_j^2}{s_j^2} + \frac{s_{j'}^2}{s_{j'}^2} - 2\rho_{jj'} = 2(1 - \rho_{jj'})
\end{aligned}$$

en donde $\rho_{jj'}$ es el coeficiente de correlación empírico entre las variables j y j' . Si este coeficiente es muy alto, aproximadamente 1, los puntos columnas están *próximos*. Si $\rho_{jj'}$ es mucho menor que 1, entonces están *alejados*.

En virtud de lo visto en AFG, las coordenadas de los puntos-columnas en el eje factorial α , son las componentes de $X'v_\alpha$ que son iguales a $u_\alpha \sqrt{\lambda_\alpha}$, en donde, en el caso de ACP normalizado, λ_α , son los autovalores de $X'X$, la matriz de correlaciones C . Recuérdese también que u_α son los autovalores de $X'X$.

Finalmente cabe observar que en el caso de la matriz considerada en ACP normalizado, la matriz de correlación $X'X$, es claro que:

$$\text{tr}(X'X) = \sum_{\alpha=1}^p \lambda_\alpha = \sum_{i,j} x_{ij}^2 = p \quad (2.5)$$

2.2.1. Ejemplo.

$$\bullet \mathbf{x}_{ij} = \mathbf{r}_{ij} - \bar{\mathbf{r}}_j$$

$$X = P \cdot R \quad p_{ij} = \begin{cases} 1 - \frac{1}{n} & i = j \\ -\frac{1}{n} & i \neq j \end{cases} \quad P_{(n \times n)} = I_{(n \times n)} - \frac{1}{n} \cdot v_{(n \times 1)} \cdot v'_{(1 \times n)} \quad v' \equiv (1, \dots, 1)$$

$$r_j \in \mathbb{R}^n ; \forall j = 1 \dots n \Rightarrow P \cdot r_j = r_j - \frac{1}{n} \cdot v \cdot (v' \cdot r_j) \quad \frac{1}{n} \cdot v' \cdot r_j = \frac{1}{n} \sum_{i=1}^n r_{ij} = \bar{r}_j \Rightarrow$$

$$P \cdot r_j = r_j - v \cdot \bar{r}_j \Rightarrow x_{ij} = r_{ij} - \bar{r}_j$$

$$A_1(1,4) \quad B_1(3,1) ; R = (A_1, B_1) = \begin{pmatrix} 1 & 3 \\ 4 & 1 \end{pmatrix} ; P = \begin{pmatrix} 1/2 & -1/2 \\ -1/2 & 1/2 \end{pmatrix} ; \bar{r}_j = (2,5,2)$$

$$\begin{pmatrix} 1/2 & -1/2 \\ -1/2 & 1/2 \end{pmatrix} \cdot \begin{pmatrix} 1 & 3 \\ 4 & 1 \end{pmatrix} = \begin{pmatrix} -3/2 & 1 \\ 3/2 & -1 \end{pmatrix} = \begin{pmatrix} 1-2,5 & 3-2 \\ 4-2,5 & 1-2 \end{pmatrix} = (A_2, B_2)$$

$$(A_2, B_2) = \begin{pmatrix} r_{11} - \bar{r}_1 & r_{12} - \bar{r}_2 \\ r_{21} - \bar{r}_1 & r_{22} - \bar{r}_2 \end{pmatrix} = \begin{pmatrix} -3/2 & 1 \\ 3/2 & -1 \end{pmatrix}$$

$$\bullet x_{ij} = \frac{r_{ij} - \bar{r}_j}{s_j \sqrt{n}}$$

$$s_j = (1,5,1) ; X = P \cdot R = \begin{pmatrix} 1/2 & -1/2 \\ -1/2 & 1/2 \end{pmatrix} \cdot \begin{pmatrix} \frac{2}{3\sqrt{2}} & \frac{3}{\sqrt{2}} \\ \frac{8}{3\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix} = \begin{pmatrix} -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{pmatrix} = (A_3, B_3)$$

$$d(A_3, B_3) = 2 \Rightarrow \text{radio} = 1$$

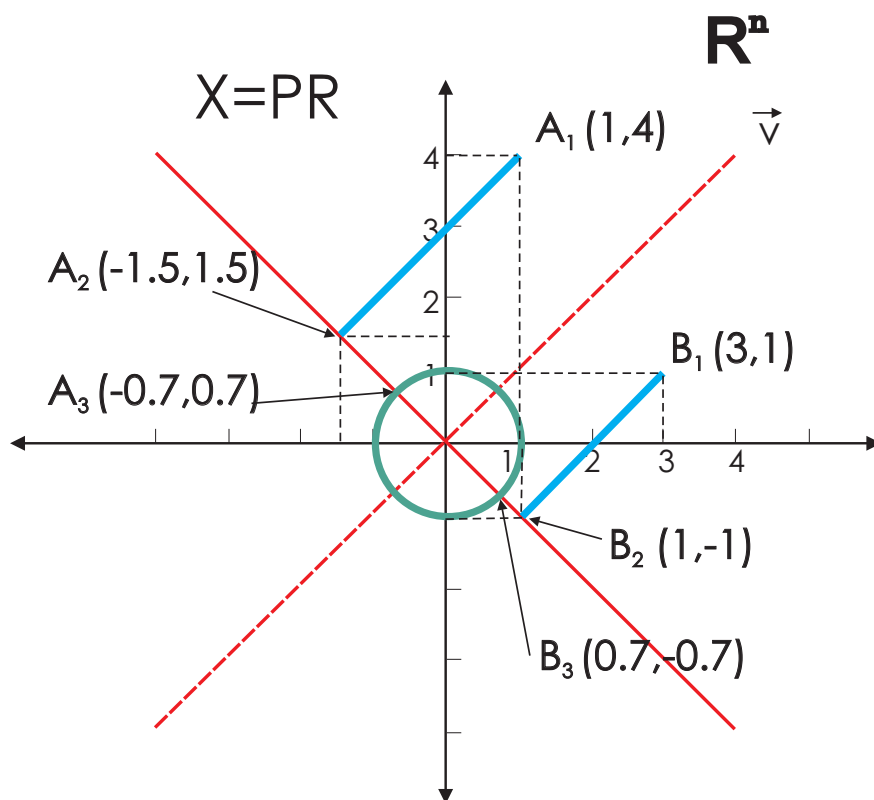


Figura 2.2: hipersfera de radio 1

2.3. Variables e individuos suplementarios en ACP.

Suele ocurrir en la práctica que, una vez realizado un ACP con una determinada matriz R de datos iniciales, se conozcan los datos en las p -variables correspondientes a nuevos individuos (Nuevos individuos que se incorporan al estudio, un grupo de individuos-control, etc.) Un caso frecuente en las aplicaciones se

produce cuando el número de individuos es muy grande —en encuestas, por ejemplo— y entonces se clasifican en grupos con arreglo a alguna o algunas características de los mismos (intervalos de ingresos familiares, por ejemplo) que sean de interés.

Incluso, a veces, interesan más estas características de los individuos que ellos mismos. En este caso, se consideran los *centros de gravedad* de las clases de individuos consideradas y estos centros se consideran *individuos nuevos* (o suplementarios) y se incorporan y presentan gráficamente como si fueran unos individuos más.

También puede darse el caso de incorporar nuevas variables cuyos valores son conocidos para los individuos estudiados, pero que en una primera etapa no se han considerado. Posiblemente, en muchos casos es así, porque no constituyen un conjunto homogéneo de características junto a las que se han seleccionado para realizar el ACP en una primera etapa.

El problema, por tanto, es cómo incorporar estos nuevos puntos-individuos (filas) y puntos-columna (variables) al ACP.

1. Consideremos los individuos suplementarios, en primer lugar. Estos constituyen un *borde* o caja matricial R_+ de la matriz R , constituida por las filas correspondientes a los individuos añadidos.

Es claro que al incrementar el número de filas de R en la forma

$$R_{n \times p} \longrightarrow \left[\begin{array}{c} R_{n \times p} \\ R_{+n' \times p} \end{array} \right]_{(n+n') \times p}$$

estas filas-individuos de R_+ han de ser comparables con las analizadas en el ACP. Es decir, hay que efectuar sobre las filas de R_+ las transformaciones que se han efectuado para realizar el ACP sobre las filas de R . Si denotamos por r_{+i} las filas suplementarias de R_+ , sería:

$$x_{+ij} = \frac{r_{+ij} - \bar{r}_j}{s_j \sqrt{n}}$$

de modo que tendríamos así definida la matriz transformada X_+ . En tal situación, entonces, las coordenadas de los individuos suplementarios sobre los α -ejes factoriales ya calculados sobre X se toman como las dadas por $X_+ u_\alpha$.

2. Análogamente procederemos con las variables suplementarias introducidas. En este caso, R se incrementa con una caja matricial R^+ de la forma

$$R_{n \times p} \longrightarrow \left[R_{n \times p} | R_{n \times p'}^+ \right]_{n \times (p+p')}$$

y en este caso, para hacer comparables las nuevas columnas (variables) en \mathbb{R}^n , con las previamente consideradas, hay que llevarlas a la esfera unidad de \mathbb{R}^n , lo que se consigue definiendo para R^+ , la matriz $X^+ = (x_{ij}^+)$, definida por

$$x_{ij}^+ = \frac{r_{ij}^+ - \bar{r}_j^+}{s_j^+ \sqrt{n}}$$

en donde \bar{r}_j^+ son las medidas de las columnas añadidas de R^+ , y s_j^+ las desviaciones típicas correspondientes.

Una vez definida X^+ , las coordenadas de estos nuevos puntos-columna, respecto del α -eje previamente calculado con X , vendrán dadas por $(X^+)' v_\alpha$.

Comentario 2.3.1 *Nótese la filosofía que se ha empleado para calcular las coordenadas $X_+ u_\alpha$ y $(X^+)' v_\alpha$ de los puntos-fila y puntos-columna suplementarios respectivamente: Se mantiene la estructura factorial obtenida a partir de la matriz inicial R sin orlar; por tanto se mantienen los vectores unitarios u_α y v_α cuyos soportes son los respectivos ejes factoriales F_α y G_α en \mathbb{R}^p y \mathbb{R}^n y, respecto de estos ejes ya obtenidos, se ubican los nuevos puntos-fila y puntos-columna. Esta metodología o filosofía de actuación respecto de los individuos suplementarios, se repetirá en otras técnicas, por ejemplo en Análisis de Correspondencias, como se verá posteriormente.*