

Introducción al Análisis Multivariado

SEMANA-5: Evaluación de la Normal Multivariada-I

Raúl Alberto Pérez

Universidad Nacional de Colombia

Escuela de Estadística

Semestre 2021-I

Oficina 43-216A

Correo: raperez1@unal.edu.co

Evaluación del Supuesto de Normalidad Multivariado

Evaluar el supuesto de normalidad multivariada es importante para facilitar los procesos de inferencia estadística.

Existen varios métodos o alternativas para evaluar el supuesto de normalidad multivariada de un vector aleatorio.

Cuando n -es grande y los métodos de evaluación utilizados se basan en el vector de medias muestrales $\bar{\mathbf{x}}$ o en ciertas distancias que involucran dicho vector de medias muestrales, el supuesto de normalidad multivariada parece no ser tan crucial.

En este caso la calidad de las inferencias realizadas, dependerá de qué tan parecida sea la distribución del vector de medias muestrales $\bar{\mathbf{x}}$ a una normal multivariada.

También es importante tener métodos para identificar cuando la distribución de un vector aleatorio se aleja de la normal multivariada, y así tener cuidado con los análisis posteriores.

Propiedad de la Normal-Multivariada:

Recordar que bajo el supuesto de normalidad multivariada de un vector aleatorio $\underline{x}_{p \times 1}$, cualquier combinación lineal de las componentes de dicho vector tiene una distribución normal univariada.

Los gráficos de contorno de la normal tri-variada son elipsoides (o hiperelipsoides) y para el caso particular de $p = 2$ (normal bi-variada) son elipses.

Preguntas Adecuadas:

Algunos pasos previos antes de evaluar la normalidad multivariada, corresponden a responder las siguientes preguntas:

1. ¿Son las distribuciones marginales del vector aleatorio, parecidas a normales univariadas?

2. ¿Es la distribución de alguna combinación lineal de las componentes de $\underline{x}_{p \times 1}$, NO parecida a una normal univariada?
3. ¿Al hacer gráficos de dispersión por pares de componentes de $\underline{x}_{p \times 1}$, presentan algunos de ellos comportamientos no elípticos?
4. ¿ Existen datos atípicos a nivel marginal o a nivel bivariado?

Métodos prácticos:

En la práctica para evaluar el supuesto de normalidad multivariado, generalmente se procede analizando la normalidad de las marginales del vector $\underline{x}_{p \times 1}$ y analizando la normalidad bivariada de pares de componentes de dicho vector.

En la práctica, no es frecuente encontrar conjuntos de datos normales en bajas dimensiones (ie, $p = 1$ o $p = 2$) y que no lo sean en altas dimensiones, pero hay que tener en cuenta que **la normalidad univariada no implica la normalidad multivariada**, ver un ejemplo en el EJERCICIO 4.8 de Johnson, un caso de normalidad univariadas que no implica Normalidad-bi-variada.

Evaluación a nivel marginal (ie. Normalidad Univariada)

Existen varios enfoques para evaluar la normalidad Uni-variada, entre ellos están los siguientes.

1. Gráficos como histogramas, cajas de bigote, etc.
2. Gráficos *qq*-plot y *NPP*.
3. Prueba de Normalidad basada en el coeficiente de correlación de los puntos del *qq*-plot.
4. Análisis de las combinaciones lineales de las variables del vector.
5. Pruebas formales de Normalidad.

A continuación se explican de manera breve cada uno de estos procedimientos.

1. Gráficos como histogramas, cajas de bigote, etc.

Generalmente se utilizan los histogramas o cajas de bigotes cuando la muestra es de tamaño moderado o grande y los diagramas de puntos en el caso de n -pequeño, para detectar alejamientos de la simetría de los datos, una cola parece ser mayor que otra.

2. Gráficos qq -plot.

Uso del gráficos cuantil-cuantil o qq -plot. Estos Son gráficos especiales que pueden ser usados para evaluar la normalidad de cada variable.

En ellos se grafican los cuantiles muestrales contra los cuantiles que se esperaría observar si las observaciones realmente provienen de una distribución normal.

Los pasos para construir un $Q - Q$ -plot son:

- (a) Ordene las observaciones originales de menor a mayor.

Sean $X_{(1)}, X_{(2)}, \dots, X_{(n)}$, las observaciones ordenadas.

Las probabilidades correspondientes a ellos son:

$$\left(1 - \frac{1}{2}\right) / n, \left(2 - \frac{1}{2}\right) / n, \dots, \left(n - \frac{1}{2}\right) / n$$

Lo anterior quiere decir, que la proporción i/n -de la muestra que está a la izquierda de $X_{(i)}$ se aproxima por: $\frac{i - \frac{1}{2}}{n}$.

- (b) Se calculan los cuantiles de la normal estándar: $q_{(1)}, q_{(2)}, \dots, q_{(n)}$, correspondientes a las probabilidades anteriores.

(c) Grafique los pares de observaciones:

$$[q_{(1)}, X_{(1)}] , [q_{(2)}, X_{(2)}] , \cdots , [q_{(n)}, X_{(n)}].$$

Si los datos proceden de una distribución normal, se espera que el gráfico anterior sea aproximadamente una línea recta.

Lo anterior significa que los pares $[q_{(i)}, X_{(i)}]$ estarán aproximadamente relacionados de forma lineal, ya que $\sigma q_{(j)} + \mu$ estará muy cerca del cuantil muestral esperado.

3. Gráficos *NPP*.

Usar el Normal-probability-plot (*NPP*), para el cual se grafican las parejas, $[m_{(i)} , x_{(i)}]$, donde $m_{(i)} = E[Z_{(i)}]$ -es el valor esperado del i -ésimo estadístico de orden en una muestra de tamaño n de una normal estándar.

Este gráfico debe dar similar al *qq*-plot, ie. una línea recta.

Estos gráficos *qq*-plot y *NPP* no son muy claro, a menos que los tamaños de muestra sean moderadamente grandes, (ie. por ejemplo $n \geq 20$), ya que pueden mostrar observaciones muy alejadas de una tendencia lineal, aún cuando se sabe que los datos provienen de una distribución normal.

Ejemplo 1 .

Considere una muestra de tamaño $n = 10$ observaciones, las cuales fueron ordenadas de menor a mayor como se muestra en la siguiente tabla.

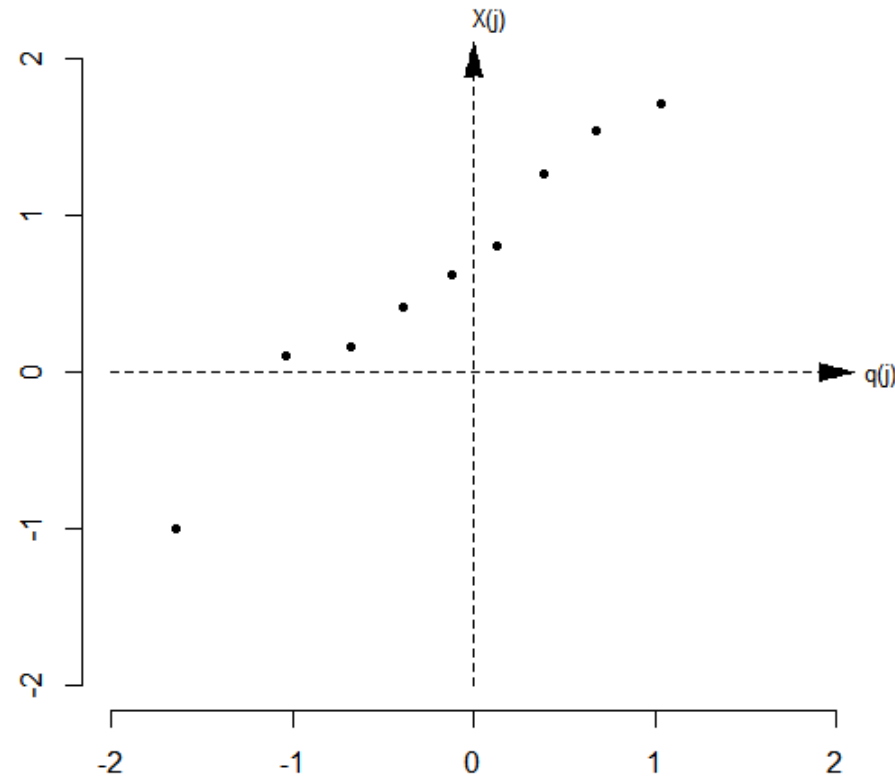
Por ejemplo, para el cálculo del cuantil de la $N(0, 1)$, para una probabilidad de 0.65 busca el cuantil que satisface:

$$P[Z < q_{(7)}] = 0.65,$$

de donde se tiene que dicho cuantil es: $q(7) = 0.385$, puesto que:

$$P[Z < 0.385] = \int_{-\infty}^{0.385} \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz = 0.65.$$

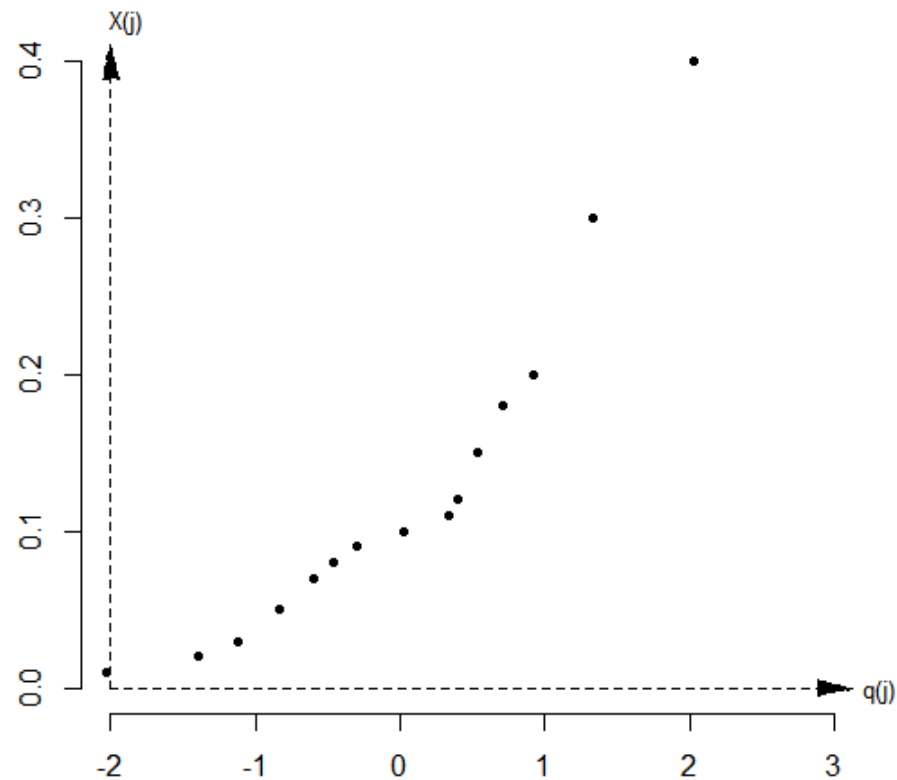
Observaciones Ordenadas $x_{(j)}$	Probabilidades $\left(j - \frac{1}{2}\right) / n$	Quantiles de la $N(0, 1)$ $q_{(j)}$
-1.00	0.05	-1.645
-0.10	0.15	-1.036
0.16	0.25	-0.674
0.41	0.35	-0.385
0.62	0.45	-0.125
0.80	0.55	0.125
1.26	0.65	0.385
1.54	0.75	0.674
1.71	0.85	1.036
2.30	0.95	1.645



La construcción del qq -plot se basa en el diagrama de dispersión de los puntos $(q_{(j)}, x_{(j)})$, $j = 1, 2, \dots, 10$, los cuales caen muy cerca de una línea recta, lo que conduce a no rechazar la hipótesis de que estos datos son de una distribución normal.

Ejemplo 2 *El departamento de control de calidad de una empresa que produce hornos micro-ondas requiere monitorear la cantidad de radiación emitida por dichos micro-ondas cuando tienen la tapa cerrada. Aleatoriamente se eligieron $n = 42$ hornos y se observó dicha cantidad de radiación emitida por ellos con la tapa cerrada.*

Horno	Radiación	Horno	Radiación	Horno	Radiación	Horno	Radiación
1	0.15	12	0.02	23	0.03	34	0.30
2	0.09	13	0.01	24	0.05	35	0.02
3	0.18	14	0.10	25	0.15	36	0.20
4	0.10	15	0.10	26	0.10	37	0.20
5	0.05	16	0.10	27	0.15	38	0.30
6	0.12	17	0.02	28	0.09	39	0.30
7	0.08	18	0.10	29	0.08	40	0.40
8	0.05	19	0.01	30	0.18	41	0.30
9	0.08	20	0.40	31	0.10	42	0.05
10	0.10	21	0.10	32	0.20		
11	0.07	22	0.05	33	0.11		



La apariencia del gráfico indica que los datos no parecen provenir de una distribución normal. Los puntos señalados con un círculo son observaciones atípicas, pues están muy lejos del resto de los datos.

Observación: Para esta muestra de datos de radiación, varias observaciones son iguales (observaciones empatadas). Cuando esto ocurre, a las observaciones con valores iguales se les asigna un mismo cuantil, el cual se obtiene usando el promedio de los cuantiles que ellas hubieran tenido si hubieran sido ligeramente distintas.

4. Prueba de Normalidad basada en el coeficiente de correlación de los puntos del *qq*-plot.

La linealidad de un gráfico *qq*-plot puede ser medida calculando el coeficiente de correlación para los puntos de dicho gráfico, el cual está dado por:

$$r_Q = \frac{\sum_{j=1}^n (x_{(j)} - \bar{x})(q_{(j)} - \bar{q})}{\sqrt{\sum_{j=1}^n (x_{(j)} - \bar{x})^2} \sqrt{\sum_{j=1}^n (q_{(j)} - \bar{q})^2}}$$

Basados en este coeficiente de correlación, se puede construir una prueba potente de normalidad (ver. Filliben, 1975; Looney y Gulledge, 1985; Shapiro y Wilk, 1965). Formalmente, se rechaza la hipótesis de normalidad a un nivel de significancia de α si $r_Q < r_Q(\alpha, n)$ donde los valores críticos $r_Q(\alpha, n)$ se encuentran en la siguiente tabla:

Valores críticos para el coeficiente de correlación del gráfico *qq*-plot para probar normalidad.

Tamaño de muestra n	Nivel de Significancia α		
	0.01	0.05	0.10
5	0.8299	0.8788	0.9032
10	0.8801	0.9198	0.9351
15	0.9126	0.9389	0.9503
20	0.9269	0.9508	0.9604
25	0.9410	0.9591	0.9665
30	0.9479	0.9652	0.9715
35	0.9538	0.9682	0.9740
40	0.9599	0.9726	0.9771
45	0.9632	0.9749	0.9792
50	0.9671	0.9768	0.9809
55	0.9695	0.9787	0.9822
60	0.9720	0.9801	0.9836
75	0.9771	0.9838	0.9866
100	0.9822	0.9873	0.9895
150	0.9879	0.9913	0.9928
200	0.9905	0.9931	0.9942
300	0.9935	0.9953	0.9960

Ejemplo 3 .

Para el primer ejemplo donde $n = 10$, el cálculo del coeficiente de correlación entre los puntos $(q_{(j)}, X_{(j)})$, $j = 1, 2, \dots, 10$, del gráfico qq-plot es:

$$r_Q = \frac{\sum_{j=1}^{10} (x_{(j)} - \bar{x})(q_{(j)} - \bar{q})}{\sqrt{\sum_{j=1}^{10} (x_{(j)} - \bar{x})^2} \sqrt{\sum_{j=1}^{10} (q_{(j)} - \bar{q})^2}} = \frac{8.584}{\sqrt{8.472} \sqrt{8.795}} = 0.994,$$

de donde, para un nivel de significancia $\alpha = 0.10$, el valor crítico de la tabla es: $r_Q(0.10, 10) = 0.9351$, luego como

$r_Q = 0.994 > 0.9351 = r_Q(0.10, 10)$, entonces no rechazamos la hipótesis de normalidad.

Observación: Para muestras grandes, las pruebas basadas en r_Q y la prueba de Shapiro Wilk (una prueba potente de normalidad) son aproximadamente las mismas.

5. Analisis de las combinaciones lineales de las variables de \underline{x}

Considere los valores propios de S , $\hat{\lambda}_1, \hat{\lambda}_2, \dots, \hat{\lambda}_p$ y sus correspondientes vectores propios $\hat{e}_1, \hat{e}_2, \dots, \hat{e}_p$. Se sugiere verificar normalidad para las combinaciones lineales dadas por:

$$\hat{e}_1^t \underline{x}_j \quad y \quad \hat{e}_p^t \underline{x}_j$$

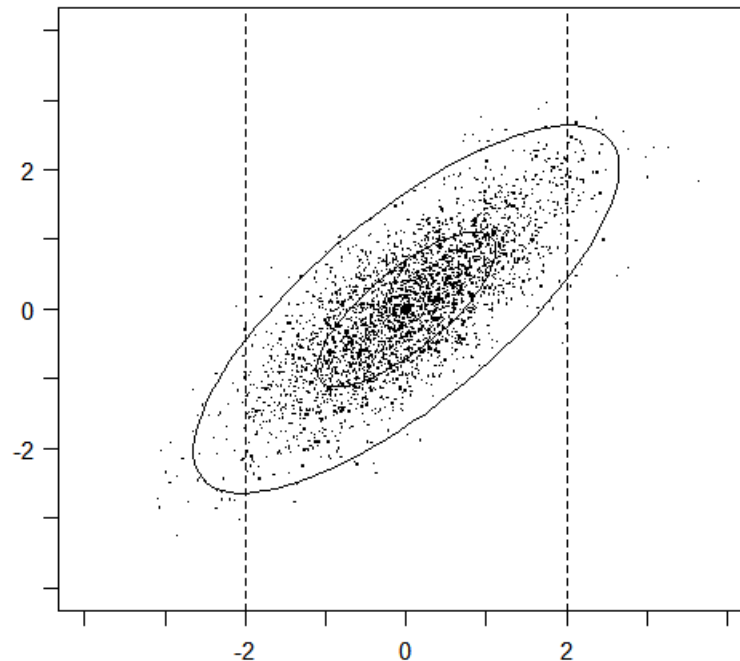
donde \hat{e}_1 y \hat{e}_p son los vectores propios correspondientes al mayor y menor valor propio de S , respectivamente.

6. Pruebas Formales.

Además de todas las ayudas anteriores, se pueden hacer pruebas formales de normalidad para el caso univariado, como lo son: Kolmogorov-Smirnov, Shapiro-Wilks, Anderson-Darling, Cramer-Von-Mises, etc.

Evaluación de la Normalidad Bi-variada

Si las observaciones fueran generadas por una distribución normal multivariada ($p > 2$), todas las distribuciones bivariadas ($p = 2$) deberían ser normales y los contornos de densidad constante deberían ser elipses. Observe el siguiente diagrama de dispersión generado por una muestra simulada de una normal bivariada.



Además, por resultado mencionado anteriormente, el conjunto de puntos bi-variados \underline{x} tales que:

$$(\underline{X} - \underline{\mu})^t \underline{\Sigma}^{-1} (\underline{X} - \underline{\mu}) \leq \chi_{\alpha}^2(2)$$

tiene una probabilidad $1 - \alpha$, ie.

$$P\left[(\underline{X} - \underline{\mu})^t \underline{\Sigma}^{-1} (\underline{X} - \underline{\mu}) \leq \chi_{\alpha}^2(2)\right] = 1 - \alpha.$$

Por ejemplo, si $\alpha = 0.5$, para muestras grandes se esperaría que alrededor del 50% de las observaciones caigan dentro de la elipse dada por:

$$(\underline{X} - \underline{\bar{x}})^t S^{-1} (\underline{X} - \underline{\bar{x}}) \leq \chi_{0.5}^2(2) = 1.39$$

Si no se cumple esto, entonces la normalidad bi-variada es sospechosa.

Ejemplo 4 Considere los pares de datos para las variables $X_1 = \text{ventas}$ y $X_2 = \text{ganancias}$ para las 10 mayores corporaciones industriales de E.U.

Compañía	X_1 -Ventas (Billones)	X_2 -Ganancias (Billones)	X_3 -Activo (Billones)
Citigroup	108.28	17.05	1484.10
General Electric	152.36	16.59	750.33
American Intl Group	95.04	10.91	766.42
Bank of America	65.45	14.14	1110.46
HSBC Group	62.97	9.52	1031.29
ExxonMobil	263.99	25.33	195.26
Royal Dutch/Shell	265.19	18.54	193.83
BP	285.06	15.73	191.11
ING Group	92.01	8.10	1175.16
Toyota Motor	165.68	11.13	211.15

Para estos datos se tiene que:

$$\bar{\mathbf{x}} = \begin{bmatrix} 155.60 \\ 14.70 \end{bmatrix}, \quad S = \begin{bmatrix} 7476.45 & 303.62 \\ 303.62 & 26.19 \end{bmatrix}$$

$$S^{-1} = \frac{1}{103623.12} \begin{bmatrix} 26.19 & -303.62 \\ -303.62 & 7476.45 \end{bmatrix}$$

Para $\alpha = 0.5$, a partir de la distribución chi-cuadrado con $p = 2$ -grados de libertad, se obtiene que: $\chi^2_{0.5}(2) = 1.39$, de donde, cualquier observación $\underline{x} = (x_1, x_2)$ que cumpla la siguiente desigualdad:

$$[x_1 - 155.60 \quad x_2 - 14.70] \left(\begin{bmatrix} 0.000253 & -0.002930 \\ -0.002930 & 0.072148 \end{bmatrix} \times 10^{-5} \right) \begin{bmatrix} x_1 - 155.60 \\ x_2 - 14.70 \end{bmatrix} \leq 1.39,$$

debe estar sobre o dentro del contorno estimado del 50% de probabilidad. De lo contrario la observación estaría por fuera.

Para las 10 observaciones observadas se tiene que sus distancias generalizadas son:

1.61, 0.30, 0.62, 1.79, 1.30, 4.38, 1.64, 3.53, 1.71 y 1.16.

Si los datos proceden de una distribución normal, se esperaría que aproximadamente el 50% de las observaciones caiga dentro o sobre el contorno estimado anterior, o dicho de otro modo, el 50% de las distancias calculadas deberían ser menores o iguales que 1.39.

Se observa que en total 4 de estas 10 distancias son menores que 1.39, lo que implica que la proporción estimada de observaciones que cumplen la desigualdad es del 40%.

La diferencia entre de esta proporción con el 50% (ie. un 10% de diferencia) proporciona una evidencia para rechazar la normalidad bivariada en estos datos. **Sin embargo, la muestra es muy pequeña para permitir obtener esta conclusión.**

Otro Método:

El procedimiento anterior es útil, pero bastante burdo. Un método más formal para evaluar la normalidad conjunta está basado en las distancias cuadráticas generalizadas, dadas por:

$$d_j^2 = (\underline{X}_j - \underline{\bar{x}})^t S^{-1} (\underline{X}_j - \underline{\bar{x}}) , \text{ para } j = 1, 2, \dots, n,$$

donde, \underline{X}_j -son las observaciones muestrales.

El siguiente procedimiento, no está limitado al caso bi-variado, por lo que puede ser usado par $p \geq 2$.

Cuando la población es normal multivariada y cuando tanto n como $n - p$ son grandes, por ejemplo, (≥ 25 o 30), cada una de las distancias al cuadrado d_j^2 -anterior, deberían comportarse como una variable aleatoria con distribución χ^2 .

Aunque, estas distancias no son independientes o distribuidas chi-cuadrados exactamente, es útil graficarlas como si lo fueran.

El gráfico resultante es llamado **Gráfico Chi-Cuadrado**, y se construye de la siguiente manera:

- Se Ordenan las distancias de menor a mayor: $d_{(1)}^2 \leq d_{(2)}^2 \leq \dots d_{(n)}^2$

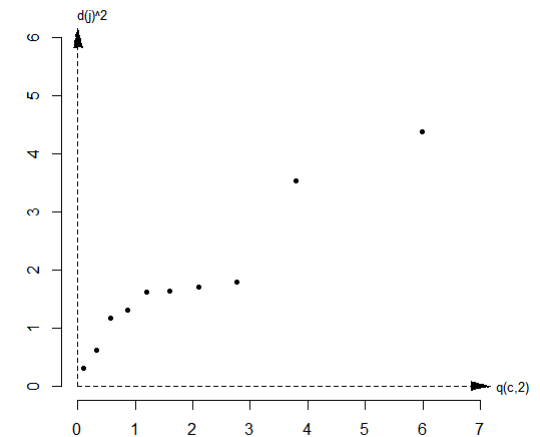
Se grafican los pares: $\left(q_{c,p} \left(\frac{j - \frac{1}{2}}{n} \right) , d_j^2 \right)$, $j = 1, 2, \dots, n$ donde, $q_{c,p} \left(\frac{j - \frac{1}{2}}{n} \right)$ es el cuantil $100 \left(j - \frac{1}{2} \right) / n$ de la distribución chi-cuadrado con p -grados de libertad, ie.

$$P \left[\chi_{(p)}^2 \leq q_{c,p} \left(\frac{j - \frac{1}{2}}{n} \right) \right] = \frac{j - \frac{1}{2}}{n}$$

Bajo normalidad, el gráfico debería mostrar un patrón lineal a través del origen y con pendiente 1. Un patrón sistemáticamente curvo sugiere falta de normalidad multivariada.

Ejemplo 5 La figura tiene las distancias ordenadas y los percentiles chi-cuadrado al igual que el gráfico chi- cuadrado, para el ejemplo anterior.

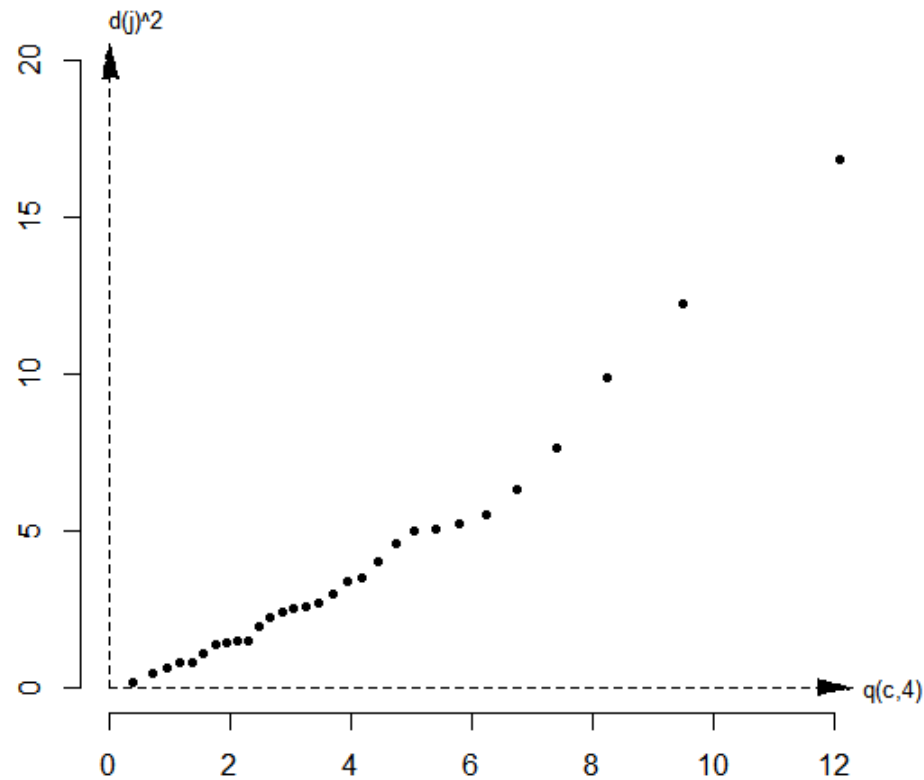
Observaciones	Distancias	Quantiles χ^2_2 de
j	Ordenadas $d^2_{(j)}$	$q_{c,2} \left(\frac{j-\frac{1}{2}}{10} \right)$
1	0.30	0.10
2	0.62	0.33
3	1.16	0.58
4	1.30	0.86
5	1.61	1.20
6	1.64	1.60
7	1.71	2.10
8	1.79	2.77
9	3.53	3.79
10	4.38	5.99



Se observa que los puntos no caen en una línea recta de pendiente 1. Por lo que no se apoya la normalidad bi-variada en estos datos. Aunque hay que tener en cuenta que n es pequeño.

Ejemplo 6 *La figura tiene las distancias ordenadas y los percentiles chi-cuadrado al igual que el gráfico chi- cuadrado, para una muestra 4-variada.*

	x_1	x_2	x_3	x_4	d_2	$q(c, 4)$		x_1	x_2	x_3	x_4	d_2	$q(c, 4)$
1	1889	1651	1561	1778	0.60	0.39	16	1954	2149	1180	1281	16.85	3.4
2	2403	2048	2087	2197	5.48	0.71	17	1325	1170	1002	1176	3.50	3.6
3	2119	1700	1815	2222	7.62	0.95	18	1419	1371	1252	1308	3.99	3.9
4	1645	1627	1110	1533	5.21	1.17	19	1828	1634	1602	1755	1.36	4.1
5	1976	1916	1614	1883	1.40	1.37	20	1725	1594	1313	1646	1.46	4.4
6	1712	1712	1439	1546	2.22	1.56	21	2276	2189	1547	2111	9.90	4.7
7	1943	1685	1271	1671	4.99	1.74	22	1899	1614	1422	1477	5.06	5.0
8	2104	1820	1717	1874	1.49	1.92	23	1633	1513	1290	1516	0.80	5.3
9	2983	2794	2412	2581	12.26	2.10	24	2061	1867	1646	2037	2.54	5.7
10	1745	1600	1384	1508	0.77	2.29	25	1856	1493	1356	1533	4.58	6.2
11	1710	1591	1518	1667	1.93	2.47	26	1727	1412	1238	1469	3.40	6.7
12	2046	1907	1627	1898	0.46	2.66	27	2168	1896	1701	1834	2.38	7.3
13	1840	1841	1595	1741	2.70	2.85	28	1655	1675	1414	1597	3.00	8.2
14	1867	1685	1493	1678	0.13	3.05	29	2326	2301	2065	2234	6.28	9.4
15	1859	1649	1389	1714	1.08	3.25	30	1490	1382	1214	1284	2.58	12.0



Se observa que los puntos no caen en una línea recta de pendiente 1. Por lo que no se apoya la normalidad bi-variada en estos datos. Aunque hay que tener en cuenta que n es pequeño.

NOTA: Consultar pruebas multivariadas de asimetría y kurtosis.

Detección de Observaciones Atípicas

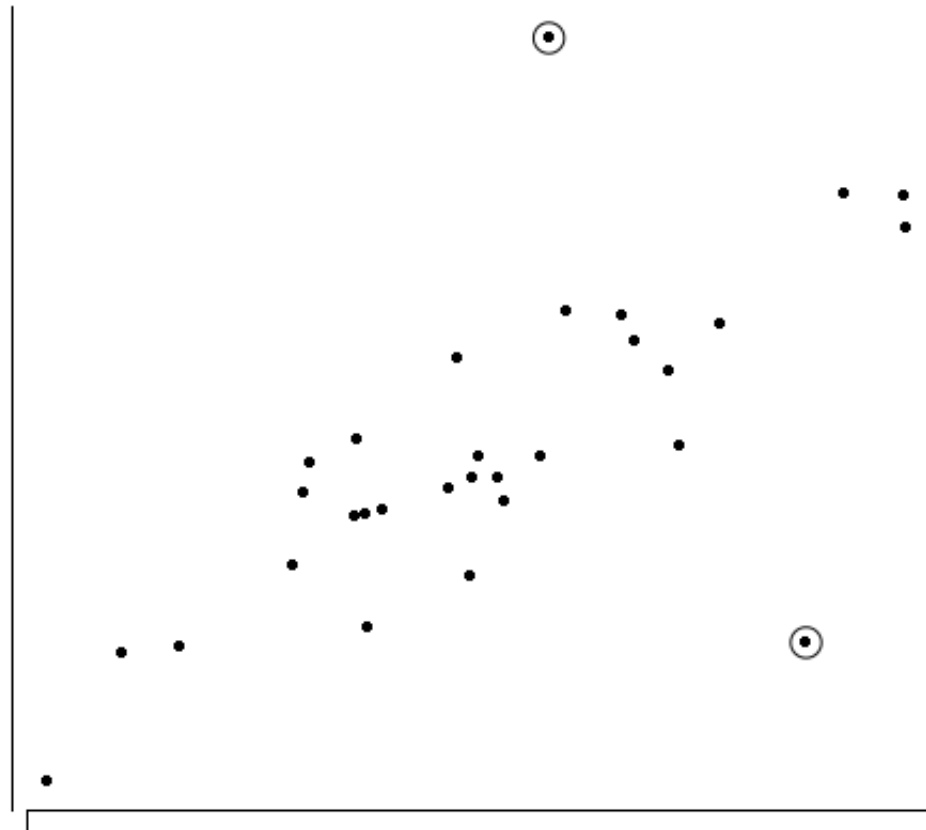
La mayoría de los conjuntos de datos contienen unas pocas observaciones inusuales que no parecen pertenecer al patrón de variabilidad seguido por las otras observaciones.

Estas observaciones son denominadas **Observaciones Atípicas** y antes de proceder a identificarlas se debe enfatizar que no todas las observaciones atípicas son números equivocados.

Elas pueden formar parte del grupo y conducir a comprender mejor el fenómeno que se estudia.

La detección de observaciones atípicas puede ser mejor realizada visualmente, es decir por medio de gráficos.

Para el caso de una variable, se pueden usar gráficos de puntos para muestras pequeñas y gráficos de cajas para muestras grandes. Para el caso de dos variables analicemos el siguiente gráfico.



Para el [Caso Bivariado](#) el diagrama de dispersión proporciona la información visual requerida para detectar datos atípicos. Sin embargo, en altas dimensiones, los datos atípicos pueden no ser detectados por gráficos uni-variados o aún diagramas de dispersión.

En estos casos se recomienda usar algunos tipos de gráficos para datos multivariados como los son: las curvas de Andrews, las gráficas de caras de Chernoff y gráficos de estrellas, etc (Consultar sección 1.4 del libro de Johnson). Estos gráficos son muy potentes para detectar casos atípicos multivariados.

[Además, en altas dimensiones un valor grande de las distancias:](#)

$$d_j^2 = (\underline{X}_j - \underline{\bar{x}})^t S^{-1} (\underline{X}_j - \underline{\bar{x}}) , \text{ para } j = 1, 2, \dots, n,$$

sugieren que la observación \underline{X}_j es inusual, aunque no la hallamos visualizado gráficamente.

Resumen: Los pasos para detectar observaciones extremas o Outliers en datos multivariados son:

a) Hacer gráficos de puntos para cada variable.

b) Hacer gráficos de dispersión para cada par de variables.

c) Calcular los valores estandarizados de cada variable dados por:

$$z_{jk} = \frac{(x_{jk} - \bar{x}_k)}{\sqrt{s_{kk}}} , \text{ para } j = 1, 2, \dots, n \text{ y } k = 1, 2, \dots, p.$$

Examinar aquellos valores estandarizados muy grandes o muy pequeños.

d) Calcular las distancias cuadradas generalizadas:

$$d_j^2 = (\underline{X}_j - \underline{\bar{x}})^t S^{-1} (\underline{X}_j - \underline{\bar{x}}) , \text{ para } j = 1, 2, \dots, n.$$

Examinar aquellas distancias de valores muy grandes.

En el gráfico chi cuadrado, son aquellos puntos más alejados del origen.

Transformaciones para acercar a la normalidad multivariada

Cuando la normalidad no es un supuesto viable, en algunos casos se pueden hacer transformaciones de los datos para acercarlos a la normalidad.

Las transformaciones son solamente reexpresiones de los datos en diferentes unidades. **Por ejemplo**, cuando un histograma de observaciones positivas muestra una gran cola derecha, una transformación de ellos tomando el logaritmo o la raíz cuadrada generalmente mejora la simetría con respecto a la media y aproxima la distribución a la normalidad.

Los tipo de transformaciones a realizar, pueden ser sugeridos por las características de los mismos datos o por consideraciones teóricas.

En el caso de transformaciones sugeridas por los mismos datos, se tiene por ejemplo que, **los datos de conteos** pueden ser más normales si se les toma la raíz cuadrada. Similarmente, para **datos de proporciones** la transformación logit y para datos de correlación la transformación de Fisher.

Resumen:

Escala Original	Escala Transformada
Conteos, y	\sqrt{y}
Proporciones \hat{p}	$Logit(\hat{p}) = \frac{1}{2}Log\left(\frac{\hat{p}}{1-\hat{p}}\right)$
Correlaciones, r	Fisher: $z(r) = \frac{1}{2}Log\left(\frac{1+r}{1-r}\right)$

En el caso de consideraciones teóricas, una familia de transformaciones para este propósito es la familia de transformaciones de potencias. Existe un método analítico conveniente para escoger una transformación de potencia dentro de dicha familia.

Familia de transformaciones de potencia de Box y Cox (1964)

$$x^{(\lambda)} = \begin{cases} \frac{x^\lambda - 1}{\lambda}, & \text{para } \lambda \neq 0 \\ \text{Ln}(x), & \text{para } \lambda = 0 \end{cases}$$

La cual es continua en λ para $x > 0$.

Dadas las observaciones x_1, x_2, \dots, x_n , la solución de Box y Cox para escoger la transformación λ adecuada (ie. el valor de λ adecuado), **es aquella que maximiza la expresión:**

$$l(\lambda) = -\frac{n}{2} \text{Ln} \left[\frac{1}{n} \sum_{j=1}^n \left(x_j^{(\lambda)} - \overline{x^{(\lambda)}} \right)^2 \right] + (\lambda - 1) \sum_{j=1}^n \text{Ln}(x_j), \text{ donde,}$$

$$\overline{x^{(\lambda)}} = \frac{1}{n} \sum_{j=1}^n x_j^{(\lambda)} = \frac{1}{n} \sum_{j=1}^n \left(\frac{x_j^{(\lambda)} - 1}{\lambda} \right),$$

es la media aritmética de las observaciones transformadas.

El proceso de maximización de $l(\lambda)$ es fácil de realizar por medio de un computador, seleccionando muchos valores diferentes para λ y calculando el respectivo valor de $l(\lambda)$.

Es útil hacer un gráfico de $l(\lambda)$ versus λ para estudiar el comportamiento en el valor del máximo $\hat{\lambda}$.

Algunos autores, recomiendan un procedimiento equivalente para encontrar λ , creando una nueva variable definida por:

$$y_j^{(\lambda)} = \frac{x_j^\lambda - 1}{\lambda \left[\left(\prod_{i=1}^n x_i \right)^{1/n} \right]^{\lambda-1}}, \quad j = 1, 2, \dots, n$$

y calculando su varianza muestral. El mínimo de la varianza ocurre en el mismo valor λ que maximiza $l(\lambda)$.

Ejemplo 7 Para un ejemplo visto anteriormente de $n = 42$ datos de la radiación de hornos de micro-ondas con la tapa cerrada, el gráfico $qq - plot$ indica que las observaciones se desvían de lo que esperaríamos si fueran normalmente distribuidas.

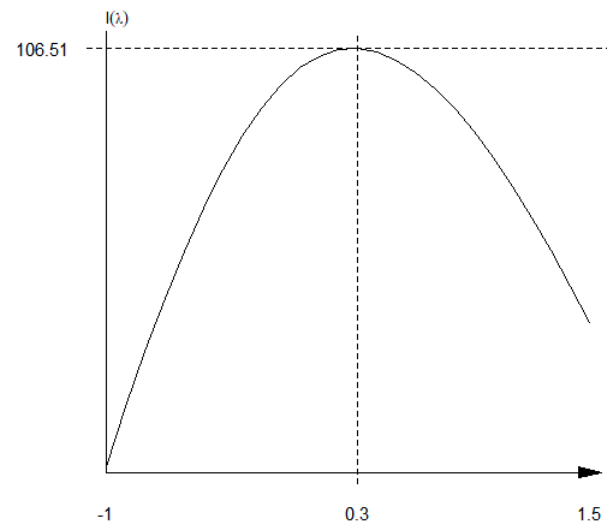
Puesto que todas las observaciones son positivas, se puede utilizar una transformación de potencia de los datos con la esperanza de acercarlos a la normalidad.

Los pares $(\lambda, l(\lambda))$, para este proceso de búsqueda se encuentran en la siguiente tabla, al igual que el gráfico de $l(\lambda)$ contra λ , donde observamos que el máximo se alcanza en $\hat{\lambda} = 0.28$, pero por conveniencia elegimos a $\hat{\lambda} = 0.3$.

	λ	$l(\lambda)$
1	-1.00	70.52
2	-0.90	75.62
3	-0.80	80.46
4	-0.70	84.94
5	-0.60	89.06
6	-0.50	92.79
7	-0.40	96.10
8	-0.30	98.97
9	-0.20	101.39

	λ	$l(\lambda)$
10	-0.10	103.35
11	0.00	104.83
12	0.10	105.84
13	0.20	106.39
14	0.30	106.51
15	0.40	106.20
16	0.50	105.50
17	0.60	104.43
18	0.70	103.03

	λ	$l(\lambda)$
19	0.80	101.33
20	0.90	99.34
21	1.00	97.10
22	1.10	94.64
23	1.20	91.96
24	1.30	89.10
25	1.40	86.07
26	1.50	82.88

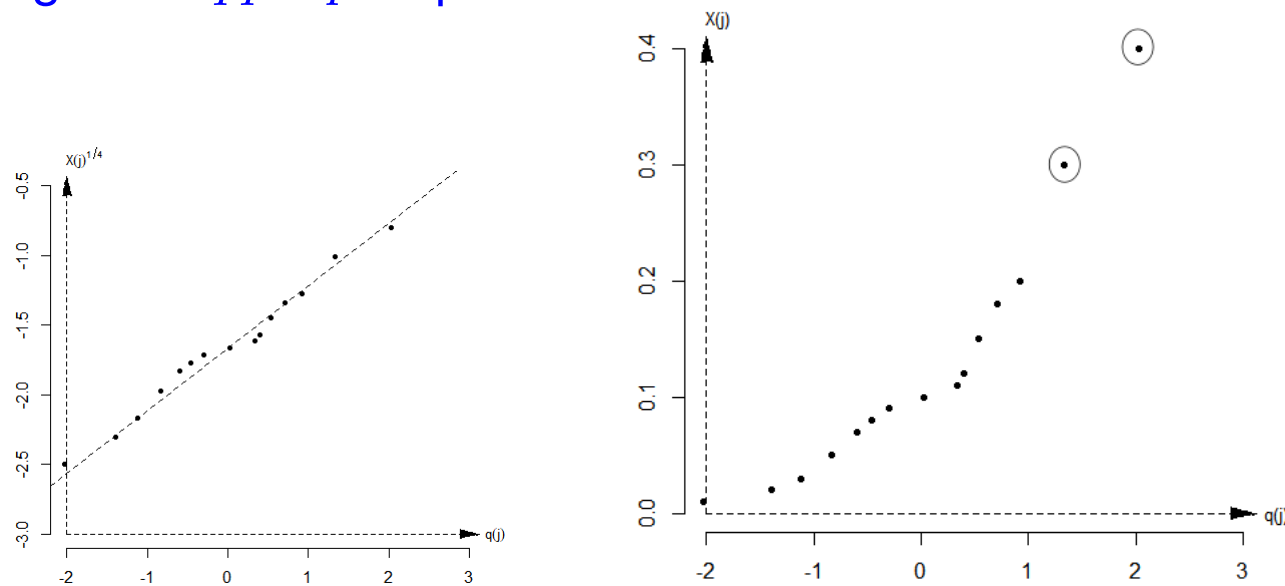


Gráfica de: $l(\lambda)$ v.s λ

Luego, los datos son transformados como:

$$x_j^{(\lambda)} = \frac{x_j^{(0.3)} - 1}{0.3}, \text{ para } j = 1, 2, \dots, 42.$$

Para verificar si los datos transformados están más cerca a la normal, se presenta su gráfico *qq - plot* para datos transformados y no-transformados.



Los pares de cuantiles caen muy cerca de una recta, lo que permite concluir que los datos: $x_j^{(\lambda)}$ son aproximadamente normal.

Transformaciones para el Caso Multivariado

Para las observaciones multivariadas se debe seleccionar una transformación para cada variable.

Sea $\lambda_1, \lambda_2, \dots, \lambda_p$ las transformaciones de potencia para las p variables consideradas.

Las transformaciones pueden ser obtenidas individualmente para cada una de las variables siguiendo el procedimiento anterior.

La j -ésima observación transformada está dada por:

$$\underline{x_j^{(\hat{\lambda})}} = \begin{bmatrix} \frac{x_{j1}^{(\hat{\lambda}_1)} - 1}{\hat{\lambda}_1} \\ \frac{x_{j2}^{(\hat{\lambda}_2)} - 1}{\hat{\lambda}_2} \\ \vdots \\ \frac{x_{jp}^{(\hat{\lambda}_p)} - 1}{\hat{\lambda}_p} \end{bmatrix}$$

donde $\hat{\lambda}_1, \hat{\lambda}_2, \dots, \hat{\lambda}_p$ son los valores que individualmente maximizan a $l(\lambda_k)$, para $k = 1, 2, \dots, p$.

El anterior procedimiento es equivalente a hacer cada distribución marginal aproximadamente normal. Aunque, la normalidad marginal de cada componente no es suficiente para garantizar que toda la distribución conjunta sea normal multivariada, frecuentemente esta condición es suficiente.

Algunos código en R