

TEST DE SHAPIRO - WILK (C7, F6, I2)

(21 / 05 / 2020, © Monfort, Dicostat2005, 2005-2020)

Le **test de SHAPIRO - WILK** est un **test de normalité** fondé sur le **rapport** de deux **estimateurs** de la **variance** qui sont des **estimateurs sans biais** sous l'hypothèse nulle, et biaisé sous l'alternative.

(i) Soit $(\mathcal{X}, \mathcal{B}, P^\xi)^{\otimes N}$ un **modèle d'échantillonnage (modèle image)** dans lequel P^ξ est la **loi de probabilité** d'une **variable parente** $\xi : \Omega \mapsto \mathbf{R}$ dont les **copies** génèrent un **échantillon indépendant équidistribué** $X = (X_1, \dots, X_N)$.

On note $X^{(\cdot)} = (X^{(1)}, \dots, X^{(N)})$ l'échantillon ordonné associé à X (cf **statistique d'ordre**) et $\mathcal{N} = \{\mathcal{N}_1(\mu, \sigma^2) : \forall (\mu, \sigma^2) \in \mathbf{R} \times \mathbf{R}_+^*\}$ la **famille** des **lois gaussiennes** (scalaires).

Pour tester la **normalité** de X , ie pour tester l'**hypothèse de base** :

$$(0) \quad H_0 : P^\xi \in \mathcal{N},$$

on peut procéder comme suit :

(a) tirage d'un **vecteur aléatoire** gaussien $U = (U_1, \dots, U_N) \sim \mathcal{N}_N(0, I_N)$, dont on déduit l'échantillon ordonné $U^{(\cdot)} = (U^{(1)}, \dots, U^{(N)})$ ainsi que le vecteur des **scores normaux** $a_N(n) = E U^{(\cdot)}$ ($n = 1, \dots, N$). La **matrice de dispersion** $\Sigma^{(\cdot)}$ de $U^{(\cdot)}$ n'est pas diagonale (cf **matrice diagonale**) car les **va** $U^{(n)}$ ne sont pas indépendantes ;

(b) **estimation**, par la **méthode des moindres carrés généralisés**, du **modèle de régression** simple :

$$(1) \quad X^{(\cdot)} = \mu \cdot e_N + \sigma \cdot U^{(\cdot)} + u,$$

dans laquelle $X^{(\cdot)}$ et $U^{(\cdot)}$ sont disposés en vecteurs colonnes. En effet, comme les $X^{(n)}$ ne sont pas indépendants, $V u$ n'est pas diagonale et l'**estimateur des mcg** de (μ, σ^2) vaut alors :

$$(2) \quad \begin{aligned} \mu_g^\wedge &= \bar{X}_N = N^{-1} \sum_{n=1}^N X_n = e_N' X / e_N' e_N \quad (\text{moyenne empirique}), \\ (\sigma^2)_g^\wedge &= \{(\mu^{(\cdot)})' (\Sigma^{(\cdot)})^{-1} \mu^{(\cdot)}\}^{-1} (\mu^{(\cdot)})' (\mu^{(\cdot)})' (\Sigma^{(\cdot)})^{-1} X^{(\cdot)}, \end{aligned}$$

où $\mu^{(\cdot)} = E U^{(\cdot)}$ et $\Sigma^{(\cdot)} = V U^{(\cdot)}$;

(c) calcul de la **statistique de S.S. SHAPIRO - M.B. WILK** (cf **statistique de test**) :

$$(3) \quad W_N^2 = T_N^2 / S_N^2,$$

dans laquelle :

$$(4) \quad \begin{aligned} S_N^2 &= X' P X, \\ T_N^2 &= \left\{ \sum_{n=1}^{\lfloor N/2 \rfloor} a_N(n) (X^{(N-n+1)} - X^{(n)})^2 \right\}. \end{aligned}$$

Le numérateur s'interprète comme l'estimateur du **coefficient de régression** de $X^{(\cdot)}$ sur $U^{(\cdot)}$, et le dénominateur comme la somme des carrés des écarts des X_n par rapport à leur **moyenne empirique** (P étant la **matrice de centrage** par rapport à cette dernière).

(ii) Le **test de S.S. SHAPIRO - M.B. WILK** est fondé sur W_N^2 et admet des **régions critiques** de la forme :

$$(5) \quad w = \{W_N^2 < q_{1-\alpha}\},$$

où $\alpha \in [0, 1[$ est le **risque de première espèce** retenu, et $q_{1-\alpha}$ le **quantile** d'ordre $1 - \alpha$ de la loi de W_N^2 .

A distance finie (ie $N \ll +\infty$), cette loi est tabulée en fonction de α et de N . En effet, on observe, par **simulation**, que $E W_N^2$ est plus grande sous l'hypothèse de **normalité** H_0 que sous une alternative non gaussienne ; de même, $V W_N^2$ est plus petite sous l'hypothèse H_0 que sous une alternative non gaussienne.

(iii) Sous l'hypothèse H_0 , on montre que :

(a) le **couple aléatoire** (S_N^2, T_N^2) est indépendant ;

(b) $W_N^2 < 1$;

(c) en posant $W_N^2 = w_N(X)$:

1) $w_N(X - \alpha) = w_N(X)$ (invariance par changement d'origine $\alpha \in \mathbf{R}^N$) ;

2) $w_N(X / \beta) = w_N(X)$ (invariance par changement d'échelle $\beta > 0$) ;

3) la loi $\mathcal{L}(W_N^2)$ de W_N^2 ne dépend pas de la **vraie valeur** $(\mu^*, (\sigma^2)^*)$ du paramètre (μ, σ^2) .

La **statistique de test** W_N^2 est souvent notée W_N (ou simplement W) et ce test est aussi appelé **W-test (de normalité)**.

(iv) La **statistique de test** de SHAPIRO-WILK se calcule aussi selon :

$$(6) \quad W_N = c_N \cdot \{S_N^2 / T_N^2\},$$

avec :

$$(4)' \quad S_N^2 = (\sigma^2)_g^{\wedge},$$

$$T_N^2 = \sum_{n=1}^N (X_n - \bar{X}_N)^2 = \|X - \bar{X}_N \cdot e_N\|^2,$$

où $c_N > 0$ désigne une **constante** de proportionnalité (cf (2)).