

Partie III

Suite du cours

"Analyse de Service"

Estimation paramétrique en survie

T v. a. de survie : f_T, S_T, F_T

C v. a. de censure : f_C, S_C, F_C

On suppose que la loi de T dépend d'un paramètre θ :

$\mathbb{P}_T = \mathbb{P}_{\theta, T}$ dépend de θ
(Loi paramétrique)

On observe T dans le cadre censuré :

on observe : $X_i = T_i \wedge C_i, \quad i=1, \dots, n$

$D = \mathbb{1}_{\{T \leq C\}}$ $T_i \ll C_i$ censure à droite

Il s'agit d'estimer le paramètre θ à partir de l'échantillon : x_1, x_2, \dots, x_n

Méthode de max. de vraisemblance

Préliminaires de calculs :

$$\text{Calcul de } \mathbb{P}(X \leq x, D=1) = \int_0^x f_T(u) S_C(u) du$$

(cf. TD)

aussi on a :

$$\mathbb{P}(X \leq x, D=0) = \int_0^x f_C(u) S_T(u) du$$

D'où les "densités" suivantes (en dérivant) :

$$\textcircled{1} \rightarrow \text{si } D=1 : f_T(x) \cdot S_C(x) \quad \text{et} \quad f_C(x) \cdot S_T(x)$$

$$\textcircled{2} \rightarrow \text{si } D=0 : f_C(x) \cdot S_T(x) \quad \text{et} \quad f_T(x) \cdot S_C(x)$$

Ainsi la densité de X_i s'écrit :

$$\left(f_T(x) S_C(x) \right)^D \cdot \left(f_C(x) S_T(x) \right)^{1-D}$$

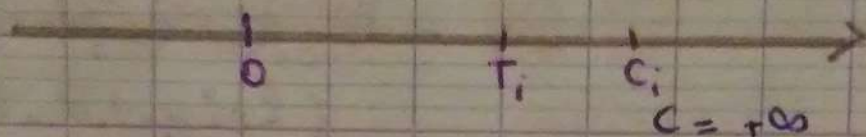
$$\partial_x \mathbb{P}(X \leq x) = \mathbb{P}(X \leq x, D=1) + \mathbb{P}(X \leq x, D=0)$$

D'où la vraisemblance de cas censuré ;

$$L_n(x_1, \dots, x_n, \theta) = \prod_{i=1}^n (f_T(x_i, \theta) S_C(x_i))^{D_i} (f_C(x_i) S_T(x_i))^{1-D_i}$$

Remarque :

Dans le cas où il n'y a pas de censure : on observe T_1, \dots, T_n



$$S_C(t) = \mathbb{P}(C > t) = 1$$

$$F_C(t) = \mathbb{P}(C \leq t) = 0$$

$$L_n(T_1, \dots, T_n, \theta) = \prod_{i=1}^n f_T(T_i, \theta)$$

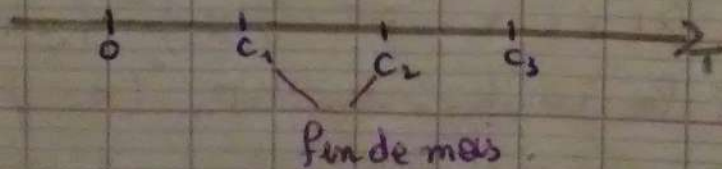
Ainsi on a l'EMV :

$$\hat{\theta}_n = \text{Arg max}_{\theta} L(x_1, \dots, x_n, \theta)$$

Exemple : $T \sim \mathcal{E}(\theta)$, C la v.a. de censure

$$C_i = c_i \quad i = 1, \dots, n$$

" constante donnée



$$f_T(t, \theta) = \theta e^{-\theta t}, \quad t > 0, \theta > 0$$

$$S_T(t, \theta) = e^{-\theta t}, \quad t > 0$$

$$f_{C_i}(t) = 1 \quad \text{si } t = c_i$$

$$= \prod_{c_i \leq t} 1 = S_{C_i}(t)$$

$$\mathbb{P}(C = c_i) = 1$$

$$S_{c_i}(t) = \prod_{L_0, c_i, t} (t)$$

D'où la vraisemblance :

$$L(x_1, x_2, \dots, x_n, \theta) = \prod_{i=1}^n (\theta e^{-\theta x_i} \prod_{(x_i > 0)} \prod_{[a, c]} (x_i))^{D_i}$$

$$= \prod_{i=1}^n (\theta e^{-\theta x_i} \prod_{[a, c]} (x_i))^{D_i} (e^{-\theta c_i})^{1-D_i}$$

$$\Rightarrow \ln L(x_1, \dots, x_n, \theta) = \log L_n(x_1, \dots, x_n, \theta)$$

$$= \sum_{i=1}^n [D_i \log(\theta e^{-\theta x_i}) \prod_{\{0 \leq x_i < c_i\}} + (1-D_i) \log e^{-\theta c_i}]$$

$$\frac{\partial \ln L}{\partial \theta} = \sum_{i=1}^n [D_i \left(\frac{1}{\theta} - x_i \right) \prod_{(0 \leq x_i < c_i)} - (1-D_i) c_i] = 0$$

$$0 = \frac{1}{\theta} \left(\sum_{i=1}^n D_i \right) - \left(\sum_{i=1}^n D_i x_i \right) - \left(\sum_{i=1}^n (1-D_i) c_i \right)$$

si $0 \leq x_i < c_i$

$$\Rightarrow \hat{\theta}_n = \frac{\sum_{i=1}^n D_i}{\sum_{i=1}^n D_i x_i + \sum_{i=1}^n (1-D_i) c_i}$$

si $0 \leq x_i < c_i$

Rem: $T \sim E(\theta)$

T_1, \dots, T_n cas non censuré

le EMV est :

$$L_n(T_1, \dots, T_n, \theta) = \prod_{i=1}^n f_T(T_i, \theta)$$

$$= \prod_{i=1}^n (\theta e^{-\theta T_i}) \prod_{(T_i > 0)} = \theta^n e^{-\theta \sum_{i=1}^n T_i}$$

si $T_i > 0$

$$L_n(\bar{T}_1, \dots, \bar{T}_n, \theta) = n \log \theta - \theta \sum_{i=1}^n T_i$$

$$\frac{\partial L_n}{\partial \theta} = \frac{n}{\theta} - \sum_{i=1}^n T_i = 0$$

$$\Rightarrow \hat{\theta}_{\text{ML}} = \frac{n}{\sum_{i=1}^n T_i} \quad (D.C. = 1)$$

§. Processus ponctuels:

. Processus de Poisson.

Soit (Y_n) v.a. i.i.d de loi $B(p)$ $0 < p < 1$

Soit (B_n) une suite de v.a. z :

$$B_0 = 0, \quad B_n = \sum_{i=1}^n Y_i$$

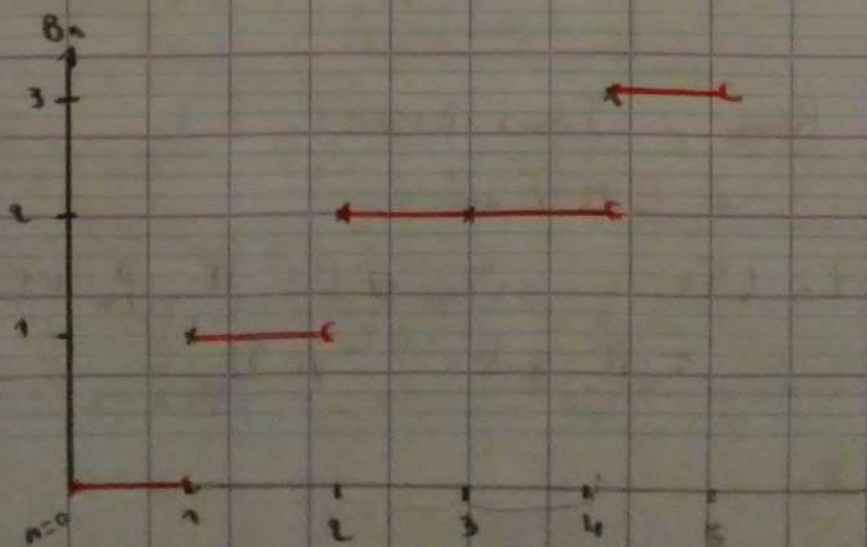
$$\text{on a: } B_n = B_{n-1} + Y_n$$

Les trajectoires de (B_n) : (on fixe la suite $(B_n(\omega))$)

n varie: 1, 2, ...

$$B_n(\omega) = \sum_{i=1}^n Y_i(\omega), \quad B_0(\omega) = 0.$$

Ainsi pour voir $B_1(\omega), B_2(\omega), \dots, B_n(\omega)$



$$Y_1(u) = 1$$

$$B_2(u) = Y_1(u) + Y_2(u), \quad B_3(u) = Y_1 + Y_2 + Y_3$$

Les projections de B_n sont ~~croissantes~~ au sens large
 $(B_n)_{n \geq 0}$ est appelé ~~le~~ processus de Bernoulli

Les instants de sauts de (B_n) : $u : 1, 2, 4, \dots$

On note les instants de sauts : S_1, S_2, \dots, S_n
 $(S_n)_{n \geq 1}$ 1^{er} saut n^{ème} saut

$$S_1(u) = 1$$

$$S_2(u) = 2$$

$$S_3(u) = 4$$

$$S_{n+1} = \inf \{ k > S_n \mid B_k \neq B_{S_n} \}$$

On pose $E_n = S_n - S_{n-1}$, $E_{n+1} = S_{n+1} - S_n$

Les écarts entre les sauts

$$E_n(u) \dots$$

On montre que (E_n) sont iid de loi $\mathcal{G}(p)$:

Loi géométrique $P(E_n = k) = p(1-p)^{k-1}$

$$\sum_{k=1}^{\infty} p(1-p)^{k-1} = p \sum_{k=1}^{\infty} (1-p)^{k-1} = p \frac{1}{1-(1-p)} = 1$$

On a aussi $S_n = E_1 + E_2 + \dots + E_n$

Car $S_n = S_{n-1} + E_n \stackrel{?}{=} S_{n-2} + E_{n-1} + E_n$

Propriétés sur la loi $g(p)$

on a l'équivalence :

$$i) \mathcal{L}(X) = g(p) \Leftrightarrow ii) \mathcal{L}(X-n | X > n) = \mathcal{L}(X)$$

caractérisation de la loi $g(p)$.

Dém : $X \sim g(p)$

$$\Rightarrow \text{on pose } G(n) = P(X > n) = \sum_{k=n+1}^{\infty} P(X=k)$$

$$= \sum_{k=n+1}^{\infty} P(1-p)^{k-1}$$

$$= \sum_{k=1}^{\infty} P(1-p)^{k-1} - \sum_{k=1}^n P(1-p)^{k-1}$$

$$Mq : P(X-n = k | X > n) = P(1-p)^{k-1} ??$$

$$P(X = n+k | X > n) = \frac{P(X = n+k, X > n)}{P(X > n)}$$

$$= \frac{P(X = n+k)}{P(X > n)} = \frac{P(1-p)^{n+k-1}}{?}$$

exercice

$$\textcircled{e} Mq \quad G(n) \cdot G(m) = G(n+m) ? \text{ exercice}$$

$$\text{dém : } G(n) = G\left(\overbrace{1}^n + \overbrace{1}^m + 1\right) = G(1) \cdot G(n-1)$$

$$= (G(1))^2 G(n-2)$$

$$= (G(1))^n$$

$$= (P(X > 1))^n$$

$$= (1 - P(X=1))^n$$

Ainsi:

$$\begin{aligned} \mathbb{P}(X = \ell) &= \mathbb{P}(X \geq \ell, X > \ell) \\ &= \mathbb{P}(X \geq \ell | X > \ell) \mathbb{P}(X > \ell) \end{aligned}$$

$$\mathbb{P}(X = \ell | X > \ell) = \mathbb{P}(X = \ell) = \dots = p(1-p)$$

$$G(n) = \mathbb{P}(X > n)$$

Autre méthode:

$$\begin{aligned} \mathbb{P}(X = \ell) &= \mathbb{P}((X \geq \ell) \setminus (X > \ell)) \\ &= \mathbb{P}(X \geq \ell) - \mathbb{P}(X > \ell) \\ &= \mathbb{P}(X > \ell - 1) - \mathbb{P}(X > \ell) \end{aligned}$$

$$= G(\ell - 1) - G(\ell) = (1-p)^{\ell-1} - (1-p)^\ell = \text{Reste à vérifier que } G(n) = 1 - p^n$$

$$\begin{aligned} \mathbb{P}(X = \ell) &= (1-p)^{\ell-1} - (1-p)^\ell = (1-p)^{\ell-1} (1 - (1-p)) \\ &= p(1-p)^{\ell-1} \end{aligned}$$

Prop 1: on a l'équivalence:

$$i) \mathcal{L}(X) = \mathcal{E}(\lambda), \lambda > 0 \Leftrightarrow ii) \mathcal{L}(X - \ell | X > \ell) = \mathcal{L}(X)$$

Prop 2: Si $X \sim g(p)$, $0 < p < 1 \forall \ell \geq 0, t \in \mathbb{R}^+$
et $Z_\ell = \frac{X}{\ell}, \ell \geq 1$

Alors si $\mathbb{E} p_\ell \xrightarrow{\ell \rightarrow \infty} \lambda > 0$ alors $Z_\ell \xrightarrow[\ell \rightarrow \infty]{\text{en loi}} Z \sim \mathcal{E}(\lambda)$

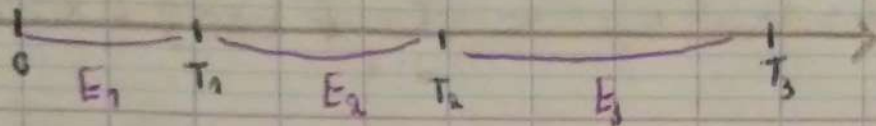
Rem: Si $Y \sim \mathcal{E}(\lambda)$ alors $[Y]$ et $Y - [Y]$ sont indép. et $1 + [Y] \sim g(e^{-1})$ exercice

(X se fait avec les fonctions génératrices)

• Processus de Poisson :

Soit (E_n) v. a. i.i.d $\rightarrow E(n)$

$(T_n)_{n \geq 0}$ v. a. n. : $T_0 = 0$, $T_n = \sum_{i=1}^n E_i$



$E_n = T_n - T_{n-1}$

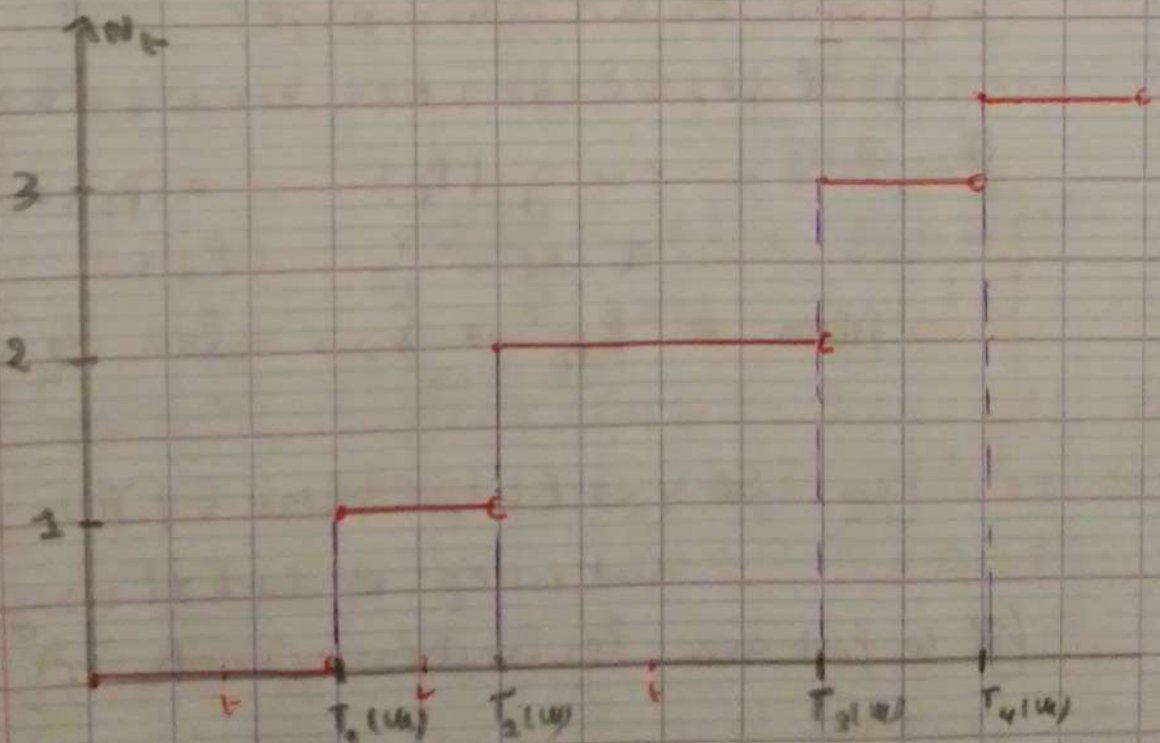
on définit N_t par :

$N_t =$ nbre de $T_i \in [0, t]$.

$N_t(\omega) = \sum_{i \geq 1} \mathbb{1}_{[0, t]}(T_i(\omega)) = \sum_{i \geq 1} \mathbb{1}_{(T_i, \infty)}(t)$

$N_t(\omega) =$ nbre de $T_i(\omega) \in [0, t]$

$(N_t)_{t \geq 0}$ est appelée le processus de Poisson (ou processus de comptage) d'intensité λ



Les trajectoires de processus de Poisson ^{sont} ~~croissantes~~ au sens large et les T_i sont les desauts de N_t .

Prop on a:

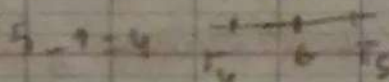
$$1) N_{T_n} = n$$

$$2) N_t = \text{Inf} \{k \geq 0 / T_{k+1} > t\}$$

$$\forall t > 0 N_t = \sum_{n \geq 0} n \mathbb{1}_{[T_n, T_{n+1}[}(t)$$

Dans cette somme tout les termes sont nuls sauf un terme si $t \in [T_n, T_{n+1}[$

$$\text{Inf} \{j-1 / T_j > t\}$$



Thm (admis): Soit $(N_t)_{t \geq 0}$ un pr. de Poisson d'intensité λ .

$(T_n)_{n \geq 1}$ les instants de sauts du ~~processus~~ de Poisson.

Alors:

1) $\mathcal{L}(T_1, T_2, \dots, T_n)$ a pour densité:

$$f(t_1, \dots, t_n) = \lambda^n e^{-\lambda t_n} \mathbb{1}_{(0 < t_1 < t_2 < \dots < t_n)}$$

2) $\mathcal{L}(T_n)$ a pour densité:

$$\forall n \geq 1 \quad f_{T_n}(t) = \frac{t^{n-1}}{(n-1)!} \lambda^n e^{-\lambda t}, t > 0$$

3) $\mathcal{L}(T_1, T_2, \dots, T_n / T_{n+1} = t)$ a pour densité:

$$f(t_1, \dots, t_n) = n! t^n \mathbb{1}_{0 < t_1 < \dots < t_n}$$

$$4) \mathcal{L}(N_t) = \mathcal{P}(\lambda t)$$

$$\mathbb{P}(N_t = \ell) = e^{-\lambda t} \frac{(\lambda t)^\ell}{\ell!}, \quad \mathbb{E}(N_t) = \lambda t$$

$$5) \mathcal{L}(T_1, \dots, T_n | N_t = n) = \text{Poi donnée au 3e.}$$

Applications : la dynamique des populations

§ Comparaison de 2 survies :

Dans certaines applications médicales, on a 2 groupes de malades G_1 et G_2 à qui on donne 2 traitements, ~~parce que~~ ~~considérés~~, on veut regarder lequel des traitements est efficace.

Si $S_2(t)$ est plus grande que $S_1(t)$ cela signifie que la survie ^{dans} G_2 est meilleur que ^{dans} G_1 c.a.d le traitement 2 est plus efficace.

on estime S_1 et S_2 par \hat{S}_{KM1} et \hat{S}_{KM2} et il s'agit de comparer ces 2 estimateurs.

Ds certains graphes on arrive à conclure tout de suite par contre ds certains les courbes se croisent.

1^{er} approche :

Il s'agit de calculer une **distance** entre les

2 survies.

Il s'agit d'un problème de test:

on veut tester $H_0: S_1 = S_2$ (égalité des 2 survies des 2 groupes),
les 2 traitements ont le même effet

contre $H_1: S_1 \neq S_2$ (effet ^{différent} des 2 traitements ~~différent~~).

Pour cela on introduit la distance suivante:

$$T_n(t) = \frac{(\hat{S}_{n1}(t) - \hat{S}_{n2}(t))^2}{V(\hat{S}_{n1}(t)) + V(\hat{S}_{n2}(t))} \quad \text{et} \quad d^*(\hat{S}_{n1}, \hat{S}_{n2}) = \sup_{t \geq 0} T_n(t)$$

n nombre d'observations

$\hat{S}_{n1}, \hat{S}_{n2}$ estimateurs de S_1 et S_2

Règle: si cette distance est "voisine" de 0 (significativement nulle)

alors on accepte H_0 .

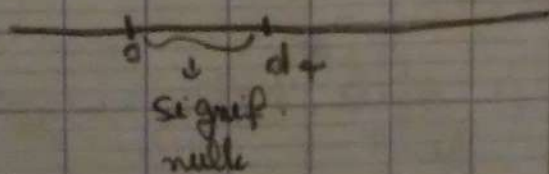
$$d^*(\hat{S}_{n1}, \hat{S}_{n2}) \approx 0 \Rightarrow \hat{S}_{n1} \approx \hat{S}_{n2}$$

pour α donné on définit une région de rejet:

$$D = \{T_n(t) > d_\alpha\}$$

Prop: sous H_0 on a $T_n(t) \rightarrow \chi^2_1$ qd n est assez grand.

$$\alpha = \mathbb{P}(T_n(t) > d_\alpha) = \mathbb{P}(\chi^2_1 > d_\alpha) \quad \text{table} \rightarrow d_\alpha$$



2^{ème} approche : test du Log-Rank

On a 2 groupes G_1 et G_2 .

Il s'agit de tester $H_0: S_1 = S_2$ (égalité de survie ds les 2 groupes)

On note : $t_1 < t_2 < \dots < t_L$ Les instants de survenue de l'év^t E ou censures des 2 groupes

Pour chaque t_k ; et $i = 1$ ou 2 ,

$n_i(t_k)$ = nombre "d'individus" à risque à T_k ds le groupe i

$d_i(t_k)$ = nbre d'événements E observés à t_k ds le groupe i

on note $d(t_k) = d_1(t_k) + d_2(t_k)$

c'est le nbre d'év E observé à T_k ds les 2 groupes

$$n(t_k) = n_1(t_k) + n_2(t_k)$$

nbre d'individus à risque à T_k ds les 2 groupes

On fait le tableau suivants (à l'instant t_k) :

	nbre d'ev E observés à t_k	Nbre d'individus à risque à t_k	nbre d'indiv en vie
G_1	$d_1(t_k)$	$n_1(t_k)$	$n_1(t_k) - d_1(t_k)$
G_2	$d_2(t_k)$	$n_2(t_k)$	$n_2(t_k) - d_2(t_k)$
Somme = Total (pour les 2 groupes)	$d(t_k)$	$n(t_k)$	$n(t_k) - d(t_k)$

On note : $O_1 =$ nbre total d'event E observés ds G_1
 $= \sum_{k | t_k \in G_1} d_1(t_k)$

$O_2 =$ nbre total d'event E observés ds G_2
 $= \sum_{k | t_k \in G_2} d_2(t_k)$

O_1 et O_2 sont observés.

on note définit :

$$E_1 := \sum_E n_1(t_k) \frac{d_1(t_k)}{n(t_k)} \quad \text{pour } G_1$$

$$E_2 := \sum_E n_2(t_k) \frac{d_2(t_k)}{n(t_k)} \quad \text{pour } G_2$$

Avec $E_{iE} := n_i(t_k) \cdot \frac{d(t_k)}{n(t_k)}$, $i = 1, 2$

$$E_{2,e} = n_2(t_e) \cdot \frac{d(t_e)}{n(t_e)} = (n_2(t_e) - n_1(t_e)) \frac{d(t_e)}{n(t_e)}$$

$$= d(t_e) - E_{1,e}$$

$$\Rightarrow E_{1,e} + E_{2,e} = d(t_e)$$

Sous H_0 : $S_1 = S_2$, on a: $O_1 + O_2 = E_1 + E_2$

$$\text{Car: } O_1 + O_2 = \sum_{e=1}^L d(t_e) = \sum_e (E_{1,e} + E_{2,e})$$

$$= \sum_e E_{1,e} + \sum_e E_{2,e} = E_1 + E_2$$

La statistique $O_1 - E_1$ (ou $O_2 - E_2$) est appelée statistique du Log Rank.

$E_{1,e}$ (ou E_2) le nombre "espéré" d'événement E observés.

On introduit la statistique suivante:

$$\frac{(O_1 - E_1)^2}{E_1} + \frac{(O_2 - E_2)^2}{E_2}$$

(distance !!)

Prop: Sous H_0 :

$$\frac{(O_1 - E_1)^2}{E_1} + \frac{(O_2 - E_2)^2}{E_2} \approx (O_1 - E_1)^2 \left(\frac{1}{E_1} + \frac{1}{E_2} \right)$$

$\rightarrow \chi^2$

D'où le test pour H_0 : $S_1 = S_2$

α donné \rightarrow test de χ^2 : da

la région de rejet est: $\mathbb{D} = \{T_n > d_\alpha\}$

Rem: Sur E_1 et E_2 : nbre "esperé" d'év^t E observés
 à t_e : $\frac{d(t_e)}{n(t_e)} \approx$ proba d'avoir "E" à t_e .
 $n(t_e) \approx p_E$

$$n_1 p_E = n_1 \frac{d(t_e)}{n(t_e)} = \underbrace{E(Z_{1i})}_{\substack{\text{d'avoir} \\ \text{l'év. E}}} \approx \beta(n_1, p_E) = E_1$$

$$E_1 = \sum_n n_1 \frac{d(t_e)}{n(t_e)}$$

On présente l'étude du log-rank sous la forme du tableau suivant

t_e	$n_1(t_e)$	$d_1(t_e)$	$n_2(t_e)$	$d_2(t_e)$	$n(t_e)$	$d(t_e)$	$n_1 \frac{d(t_e)}{n(t_e)}$	$n_2 \frac{d(t_e)}{n(t_e)}$
		$O_1 = \sum_1 d_1(t_e)$		$O_2 = \sum_2 d_2(t_e)$		$E_1 = \sum_1$		$E_2 = \sum_2$