

COURS de M1 (Stat - Proba):

Analyse de survie

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n z_i \Rightarrow N(0, F(t) \cdot S(t))$$

$$\text{or } \frac{1}{\sqrt{n}} \left( \sum_{i=1}^n \mathbb{1}_{(T_i \leq t)} - F(t) \right) \Rightarrow N(0, F(t) S(t))$$

$$\sqrt{n} (F_n(t) - F(t)) \Rightarrow N(0, F(t) S(t))$$

Sym'et

$$\left\{ \begin{aligned} &= -\sqrt{n} (S_n(t) - S(t)) \Rightarrow N(0, F(t) S(t)) \\ &= \sqrt{n} (S_n(t) - S(t)) \Rightarrow N(0, F(t) S(t)) \end{aligned} \right.$$

### §. Estimateur de Kaplan-Meier de s:

dans le paragraphe précédent, les observations  $T_1, T_2, \dots, T_n$  sont "observées" échantillon complet.

Il existe des échantillon où les  $T_n$  ne sont pas tous "observés" échantillon incomplet ou données incomplètes (observations censurées).

Soit  $T_1, T_2, T_3^+, T_4, T_5^+, T_6, T_7^+, T_{20}^+ \dots (*)$

où  $T_i^+ = T_i \wedge c_i$  où  $c_i$  v.a. de censure  
 $= X_i$

Comment estimer  $S$  ?

$$T_1, \dots, T_n \quad S_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{(T_i > t)} \xrightarrow{n \rightarrow \infty} S(t)$$

(\*)  $\rightarrow S_n(t) = \underbrace{S_{1/n}(t)}_{\substack{\text{contient les } T_i \text{ observés} \\ \text{4 termes}}} + \underbrace{S_{2/n}(t)}_{\substack{\text{sur les } T_i \text{ non obs} \\ \text{16 termes}}}$



$$S_n^+(t) = \frac{1}{20} (\underbrace{T_1 + T_2 + T_4 + T_6}_{S_1(t)}) + \frac{1}{20} (\underbrace{T_3 + T_5 + T_7}_{S_2(t) \neq S(t)})$$

la présence de donnée incomplète modifie la limite  $S(t)$ . (voir exercice)

d'où la difficulté d'estimer  $S(t)$  de la présence de donnée incomplète.

Kaplan-Meier ont proposé un estimateur de  $S$ .

On regarde d'abord 1<sup>er</sup> cas :

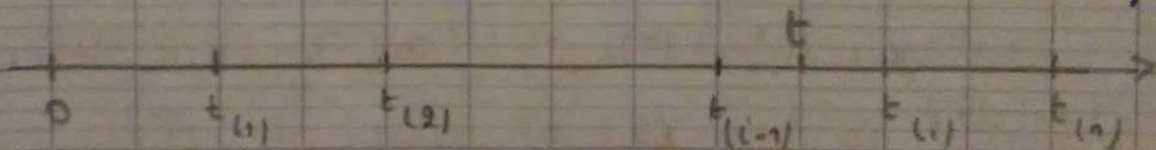
1<sup>er</sup> cas données complètes :

ech :  $T_1, T_2, \dots, T_n$  de  $T$

les réalisations :  $T_1(u), T_2(u), \dots, T_n(u)$

$$t_i = T_i(u)$$

$t_{(1)} < t_{(2)} < \dots < t_{(n)}$  instants de réalisation de "E" (événement d'intérêt).  
 $P(t_i = t_j) = 0$



pb: estimer  $S(t)$ : Soit  $t_{(i-1)} < t < t_{(i)}$  i fixe

$$S(t) = P(T > t) = IP(T > t | T > t_{(i-1)}) \cdot P(T > t_{(i-1)})$$

or :  $\otimes IP(T > t | T > t_{(i-1)}) \rightarrow$  proba que E ne se réalise pas ds

$$] t_{(i-1)}, t ]$$

$$P(T > t | T > t_{(i-1)}) = P(T > t \cap T > t_{(i-1)})$$



Si je suis "en vie" à l'instant  $t_{(i-1)}$ , c'est la proba que je reste en vie à l'instant  $t$ .

On continue  $IP(T > t_{(i-1)})$

$$IP(T > t_{(i-1)}) = IP(T > t_{(i-1)} | T > t_{(i-2)}) IP(T > t_{(i-2)})$$

donc :

$$S(t) = IP(T > t | T > t_{(i-1)}) IP(T > t_{(i-1)} | T > t_{(i-2)}) \dots IP(T > t_{(i-2)})$$

on continue la règle

$$q_j = IP(T > t_{(j)} | T > t_{(j-1)}) = q_{t_{(j)}}$$

Ainsi :

$$S(t) = IP(T > t | T > t_{(i-1)}) q_{i-1} \cdot q_{i-2} \dots q_2 IP(T > t_{(1)})$$

$$t_{(0)} = 0$$

$$q_1 = IP(T > t_{(1)} | T > 0)$$

$$IP(T > t_1) = IP(T > t_{(1)} | T > t_{(0)=0}) \cdot P(T > t_{(0)}) \\ = q_1 \cdot \underbrace{IP(T > 0)}_{=1} = q_1$$

Donc

$$S(t) = IP(T > t | T > t_{(i-1)}) \cdot q_{i-1} \cdot q_{i-2} \dots q_1$$

on note  $q_t := IP(T > t | T > t_{(i-1)})$ ,  $q_j = q_{t_{(j)}}$

alors :  $S(t) = q_t \cdot q_{i-1} \dots q_1$

$$= q_t \cdot q_{t_{(i-1)}}$$

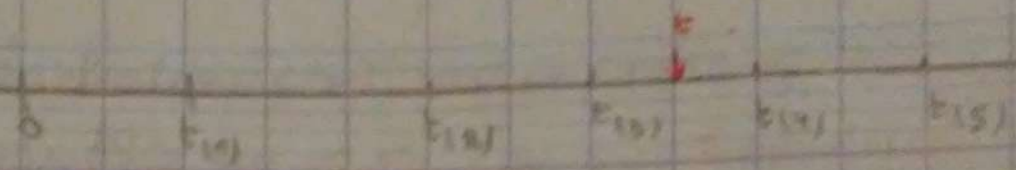
$$q_{t_{(i-1)}} = \prod_{u=1}^{i-1} q_u$$

instant

$$= S(t)$$

$$\prod_{u=1}^{i-1} q_u = q_t$$



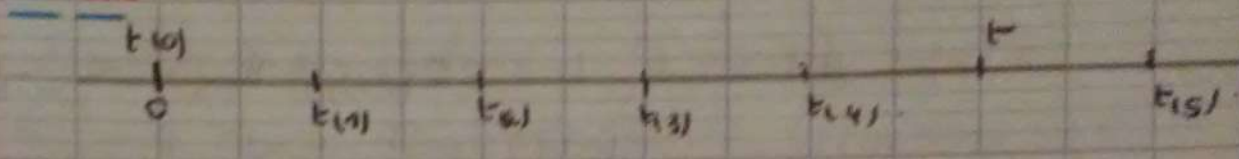


$$t_{(3)} < t < t_{(4)}$$

$$S(t) = q_0 \cdot q_{t_{(1)}} \cdot q_{t_{(2)}} \cdot q_{t_{(3)}}$$

Rem :

q pas l'événement E  
à l'instant t



$$S(t) = q_t \cdot q_4 \cdot q_3 \cdot q_2 \cdot q_1$$

c'est la survie à l'instant t : c'est ne pas avoir l'événement E à l'instant t, on a pas eu l'événement dans  $[t, t_{(4)}]$ ,  $[t_{(4)}, t_{(3)}]$ ,  $[t_{(3)}, 0]$ .

probleme : Comment estimer  $q_{t_{(j)}}$  ? pour estimer  $S(t)$ .

Estimation de  $q_{t_{(j)}}$  =  $q_j$

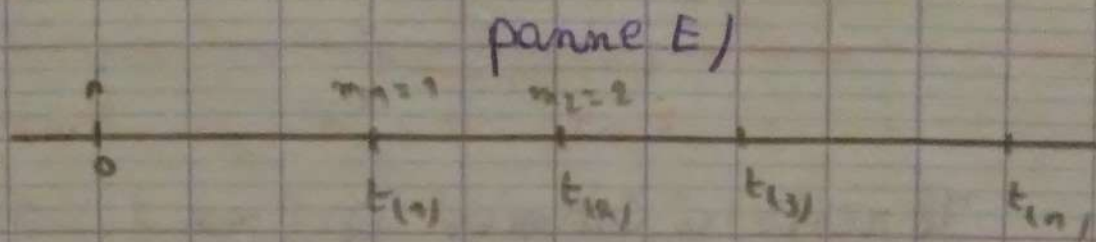
posons  $P_j = 1 - q_j$  = proba d'avoir l'événement E dans  $]t_{(j-1)}, t_{(j)}]$

on va estimer  $P_j$  par les fréquences de l'événement E  $P_j = \frac{\text{cas fav.}}{\text{cas poss.}}$

Soit  $N_j$  = nombre "d'individus" en vie à l'instant  $t_{(j-1)}$



$m_j$  : nombre d'événement E à l'instant  $t_{(j)}$   
Rem:  $n$  machines en fiabilité ( $m_j \geq 1$ ) (On observe la



$$t_{(0)} = 0 \quad P(w / x_i = x_j) = 0$$

$$n_1 \rightsquigarrow t_{(0)} = 0 \quad P(n_0) = 0$$

$$n_1 = n \quad \text{n'est pas vide}$$

$$n_2 \rightsquigarrow t_{(1)} \quad n_2 = n - 1$$

$$n_3 \rightsquigarrow t_{(2)} \quad n_3 = n - 3$$

Donc :

un estimateur de  $p_j$

$$\hat{p}_j = \frac{\text{cas favorable ds } ] t_{(j-1)}, t_{(j)} ]}{\text{cas possible}} = \frac{m_j}{n_j}$$

nbre des individus en vie à  $t_{(j-1)}$

donc  $\hat{q}_j = 1 - p_j = 1 - \frac{m_j}{n_j} \quad j=1, \dots, n.$

Ainsi l'estimateur de K. M.

on a :  $S(t_{(j)}) = \prod_{t_{(i)} \leq t_{(j)}} \hat{q}_i = q_1 q_2 \dots q_j$

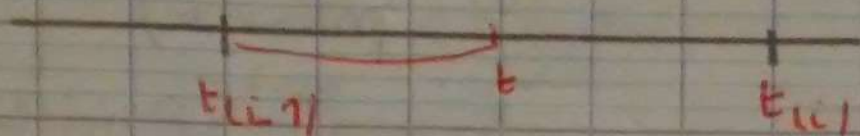
et un estimateur de  $S(t)$  est :

$$S_{KM}(t_{(j)}) = \prod_{t_{(i)} \leq t_{(j)}} \hat{q}_i = \prod_{1 \leq i \leq j} \left( 1 - \frac{m_i}{n_i} \right)$$

appelé estimateur de Kaplan Meier de  $S(t_{(j)})$



$S_{KM}(t_{(i)})$  estimateur de  $S(t_{(i)})$   
 pour  $S(t)$  où  $t_{(i-1)} < t < t_{(i)}$ .



Il reste à estimer  $q_t \rightsquigarrow$

$$q_t = \mathbb{P}(T > t \mid T > t_{(i-1)})$$

$$p_F = 1 - q_t = \text{probab d'avoir } E \text{ ds } ]t_{(i-1)}, t]$$

$$q_t = 1$$

$$S_{KM}(t) = \prod_{t_{(j)} \leq t} \hat{q}_j = \prod_{t_{(j)} \leq t} \left(1 - \frac{m_j}{n_j}\right)$$

$S_{KM}(t)$  est constant sur chaque interval  $]t_{(i-1)}, t_{(i)}]$   $i = 1, \dots, n$  et décroissante

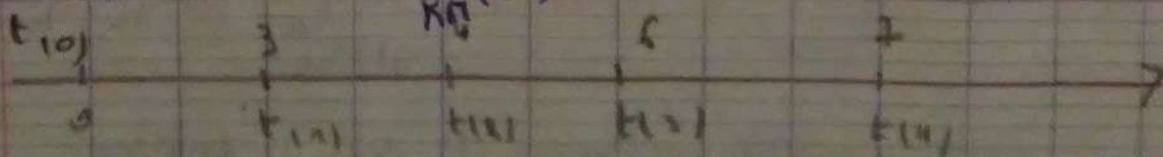
$$q_{t_0} = q_0 = 1$$

$$q_{t_1} = q_{t_0} = 1$$

Ex: on observe 5 machines et leurs durées de fonctionnement jusqu'à la panne

3, 4, 5, 6, 7 (mois)

calculer  $S_{KM}(t)$   $t > 0$



$$S_{KM}(t_{(i)}) = \prod_{d \leq i} \left(1 - \frac{m_d}{n_d}\right)$$



$m_j$ : nbr d'éven. "E" à l'instant  $t_{(j)}$   
 $n_j$ : nbr d'individus en vie à l'instant  $t_{(j-1)}$

$$m_1 = 1, m_2 = 1, m_3 = 2, m_4 = 1, m_5 = 1$$

$$n_1 = 5, n_2 = 4, n_3 = 3, n_4 = n_5 = 1, n_6 = 1$$

si  $t \in ]0, 3[$ :  $S_{KM}(t) = 1$

$$S_{KM}(3) = \left(1 - \frac{m_1}{n_1}\right) = 1 - \frac{1}{5} = \frac{4}{5} = q_1$$

si  $t \in ]3, 4[$

$$S_{KM}(t) = \frac{4}{5} = q_1$$

$$S_{KM}(4) = q_1 \times q_2 = \left(1 - \frac{m_2}{n_2}\right) \left(\frac{4}{5}\right) = \frac{3}{4} \cdot \frac{4}{5} = \frac{3}{5}$$

$$n_0 = 5; n_1 = 5, n_2 = 4, n_3 = 3, n_4 = 2, n_5 = 1$$

$$m_0 = 0, m_1 = 1, m_2 = 1, m_3 = 1, m_4 = 1, m_5 = 1$$

$$S_{KM}(t_{(0)}) = \left(1 - \frac{m_0}{n_0}\right) = 1 = S(0)$$

$$S_{KM}(t_{(1)}) = \left(1 - \frac{m_0}{n_0}\right) \left(1 - \frac{m_1}{n_1}\right) = \frac{4}{5} = S(3)$$

$$S_{KM}(t_{(2)}) = \frac{4}{5} \left(1 - \frac{m_2}{n_2}\right) = \frac{3}{5} = S(4)$$

$$S_{KM}(t_{(3)}) = \frac{3}{5} \left(1 - \frac{m_3}{n_3}\right) = \frac{2}{5} = S(5)$$

$$S_{KM}(t_{(4)}) = \frac{2}{5} = S(6)$$

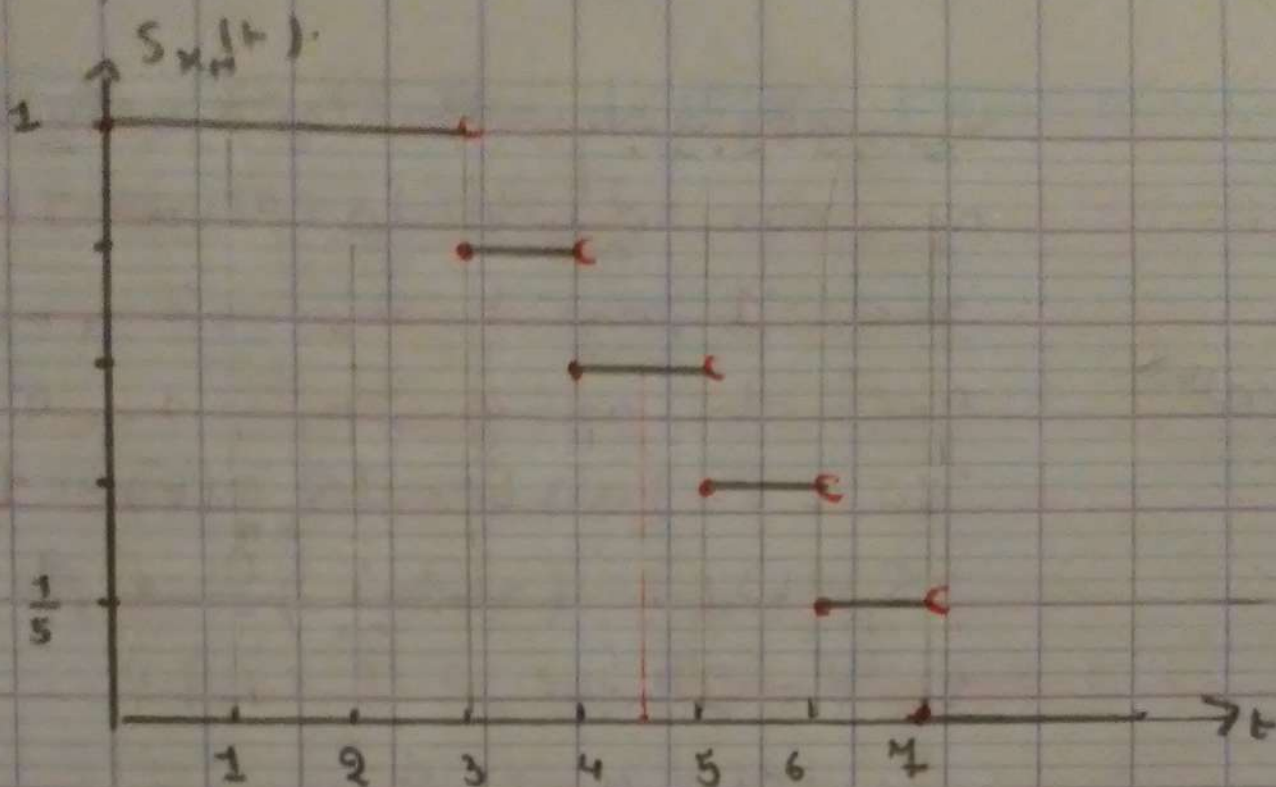
$$S_{KM}(t_{(5)}) = \frac{2}{5} \left(1 - \frac{m_5}{n_5}\right) = \frac{1}{5} (1 - 1) = 0 = S(7)$$

graphe:  $S(4, 5) = \frac{3}{5} = 0,6 = 60\%$

$t \in ]0, 3[$   $S(0) = 1$

Faux





$$S(4,5) = \frac{3}{5} = 60\% = P(T > 4,5)$$

$$S(t) = P(T > t), \quad S(5,5) = \frac{2}{5} = 40\%$$

Exemple 2)

$$\underbrace{3}_{t_{(1)}}, \underbrace{4}_{t_{(2)}}, \underbrace{6, 6}_{t_{(3)}}, \underbrace{7}_{t_{(4)}}$$

$$n_0 = 5, \quad n_1 = 5, \quad n_2 = 4, \quad n_3 = 3, \quad n_4 = 1$$

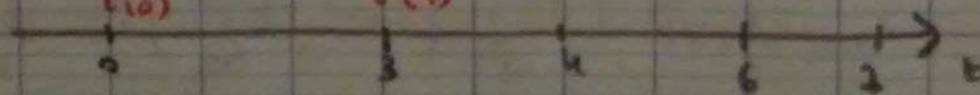
$$m_0 = 0, \quad m_1 = 1, \quad m_2 = 2, \quad m_3 = 2, \quad m_4 = 1$$

$$S_{KM}(t_{(0)}) = \frac{4}{5} = S(3)$$

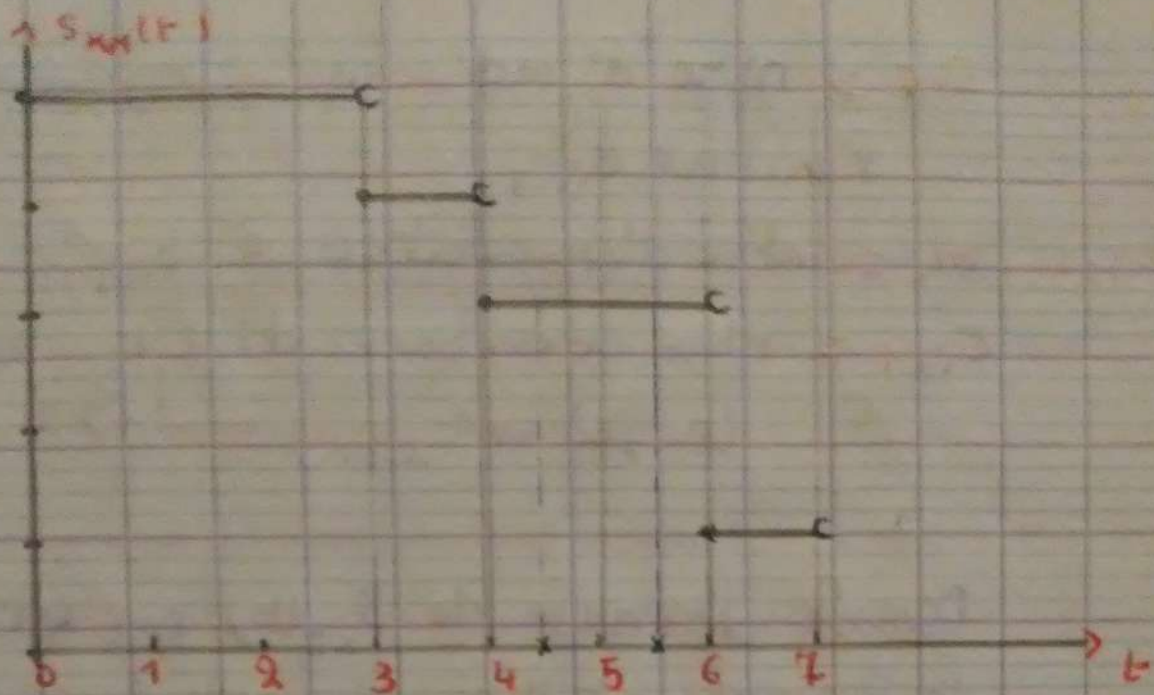
$$S_{KM}(t_{(2)}) = \frac{4}{5} \left(1 - \frac{1}{4}\right) = \frac{3}{5} = S(4)$$

$$S_{KM}(t_{(3)}) = \frac{3}{5} \left(1 - \frac{2}{3}\right) = \frac{1}{5} = S(6)$$

$$S_{KM}(t_{(4)}) = \frac{1}{5} (1 - 1) = 0 = S(7)$$







$$S(4,5) = \frac{3}{5} = 60\%$$

$$S(5,5) = \frac{3}{5} = 60\%$$

La survie est grande par rapport au 1<sup>er</sup> exemple car on a pas un év<sup>nt</sup> en 5 par contre en exemple il y avait un événement

ex:  $S_{KM}(8) = 0$ , à 8: plus "personnes" en survie.

2<sup>ème</sup> cas: données incomplètes (avec des censures):

On suppose qu'on a un échantillon d'observation avec existence de donnée censure.

On range par ordre croissant les instants  $t_{(i)}$

( $t_{(i)}$  = instant pour l'év<sup>nt</sup> E ou une censure)

les données censurées sont signalées par (+), 7<sup>+</sup>

ex: (3, 5, 6, 6, 7<sup>+</sup>, 8, 10)

à 7 on a observé une censure.  
on a observé E à ces instants



$n_i$  = nbre "d'individus" à risque (en vie) à l'instant  $t_{(i-1)}$ .

$m_i$  = nbre d'événements  $E$  à  $t_{(i)}$ .

$c_{i-1}$  = nbre de censures ds  $]t_{(i-1)}, t_{(i)}[$   
(en fait en  $t_{(i-1)}$ )  $i \geq 2$

Dans la formule de  $S_{K\pi}(t)$  qui estime  $S(t) = P(T > t)$   
la formule utilisée pour calculer  $S(t)$  dans  
le 1<sup>er</sup> cas reste la même.

$$\text{on arrive à } S_{K\pi}(t) = \prod_{t_{(j)} \leq t} \left(1 - \frac{m_j}{n_j}\right) = \prod_{t_{(j)} \leq t} q_j$$

la différence avec le 1<sup>er</sup> cas c'est l'estimation  $\hat{q}_j$

$$S(t) = q_t \quad \prod q_j = \prod q_j$$

$$q_j = 1 - p_j \rightarrow q_j = 1 - p_j$$

$$p_j = \text{proba } E \text{ ds } ]t_{(j-1)}, t_{(j)}[$$

$$\rightarrow p_j = \frac{m_j}{n_j}$$

on a :

$$t_{(0)} = 0, m_0 = 0$$

$$c_0 = \text{nbre de censure } ]t_{(0)}, t_{(1)}[ \text{ (convention)}$$

1<sup>er</sup> cas :

$$\text{on a : } n_i = n_{i-1} - m_{i-1}$$

$$n_2 = n_1 - m_1$$



2<sup>e</sup>me cas:

on a :  $t_{(0)} = 0, m_0 = 0$

$c_0 =$  nombre de censure

$] t_{(0)}, t_{(1)} [$

on a :  $n_i = n_{i-1} - m_{i-1} - c_{i-1}$

$c_0 = 0$

$n_{i-2}$

$n_i$

$t_{(i-2)}$

$t_{(i-1)}$

$t_{(i)}$

$$n_i = n_{i-1} - m_{i-1} - c_{i-1}$$

$$n_1 = n_0 - \frac{m_0}{0} - \frac{c_0}{0} \quad n_1 = n$$

$$n_2 = n_1 - m_1 - c_1$$

on a :

$$n_i = n - \sum_{j=1}^{i-1} m_j - \sum_{j=1}^{i-1} c_j$$

(car :

$$n_i = n_{i-1} - m_{i-1} - c_{i-1}$$

$$= (n_{i-2} - m_{i-2} - c_{i-2}) - m_{i-1} - c_{i-1}$$

$$= n_{i-2} - (m_{i-1} + m_{i-2}) - (c_{i-1} + c_{i-2})$$

$$\vdots$$
$$= n_0$$

L'estimateur de KM s'écrit :

$$S_{KM}(t) = \prod_{t_{(i)} \leq t} \left(1 - \frac{m_i}{n_i}\right)^{D_i} \quad (**)$$



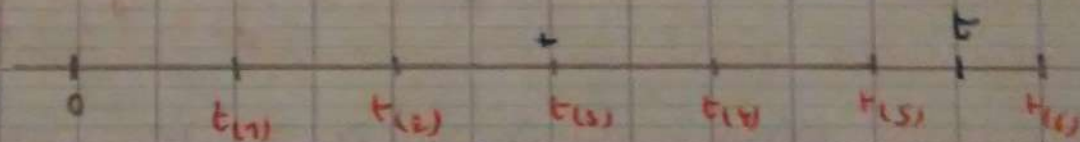
où  $D_i = \begin{cases} 1 & \text{si l'événement } E \text{ est observé en } t_{(i)} \\ 0 & \text{si il y a une censure en } t_{(i)} \end{cases}$

**Rq:** dans la formule de  $S_{KM}$  le facteur

$$\left(1 - \frac{m_i}{n_i}\right)^{D_i} = \begin{cases} 1 & \text{si } D_i = 0 \\ 1 - \frac{m_i}{n_i} & \text{si } D_i = 1 \end{cases}$$

ce qui signifie qu'il y a pas d'effet dans le calcul de produit en un point où il y a une censure mais il y a un changement dans le calcul de  $n_i$ .

$$S_{KM}(t) = \left(1 - \frac{m_1}{n_1}\right)^{D_1} \left(1 - \frac{m_2}{n_2}\right)^{D_2} \left(1 - \frac{m_3}{n_3}\right)^{D_3} \dots \left(1 - \frac{m_5}{n_5}\right)^{D_5}$$



$$S_{KM}(t) = \left(1 - \frac{m_1}{n_1}\right) \left(1 - \frac{m_2}{n_2}\right) \cdot 1 \cdot \left(1 - \frac{m_4}{n_4}\right) \left(1 - \frac{m_5}{n_5}\right)$$

pour la formule si il y'a pas de censure (on observe E d tous les  $t_{(i)}$ ) donc on retrouve la formule du 1<sup>er</sup> cas.

si il y'a des censures : le facteur  $q_j$  de

$S_{KM}$  est égale à 1 (car  $p_j = 1 - q_j = 0$ )  
 $p_j = \text{proba d'avoir } E \text{ ds } ]t_{(j-1)}, t_{(j)}$



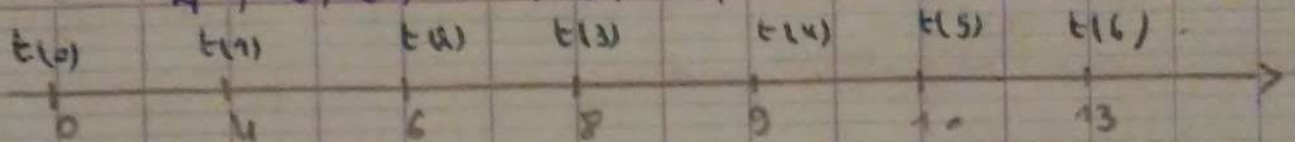
cette proba, il vaut 0.

$$P_3 \rightarrow ]t_{(2)}, t_{(3)}] \quad P_3 = 0 \quad q_3 = 1$$

exp: on observe le fonctionnement d'un nombre  $n$  de machine jusqu'à la panne.

$n = 7$  et on a les données suivantes.

4, 6, 6, 8, 9, 10, 13.



$$n_0 = 7 \quad m_0 = 0 \quad C_0 = 0$$

$$n_1 = 7 \quad m_1 = 1 \quad C_1 = 0$$

$$n_2 = 6 \quad m_2 = 2 \quad C_2 = 0$$

$$n_3 = 4 \quad m_3 = 0 \quad C_3 = 1$$

$$n_4 = 3 \quad m_4 = 1 \quad C_4 = 0$$

$$n_5 = 2 \quad m_5 = 1 \quad C_5 = 0$$

$$n_6 = 1 \quad m_6 = 1 \quad C_6 = 0$$

$$S_{KH}(t) = 1, \quad t \in [0, 4[$$

$$S_{KH}(t) = 1 \cdot \left(1 - \frac{m_1}{n_1}\right)^1 = 1 - \frac{1}{7} = \frac{6}{7} \quad t \in [4, 6[$$

$$S_{KH}(t) = \frac{6}{7} \left(1 - \frac{m_2}{n_2}\right)^2 = \frac{4}{7} \quad t \in [6, 8[$$

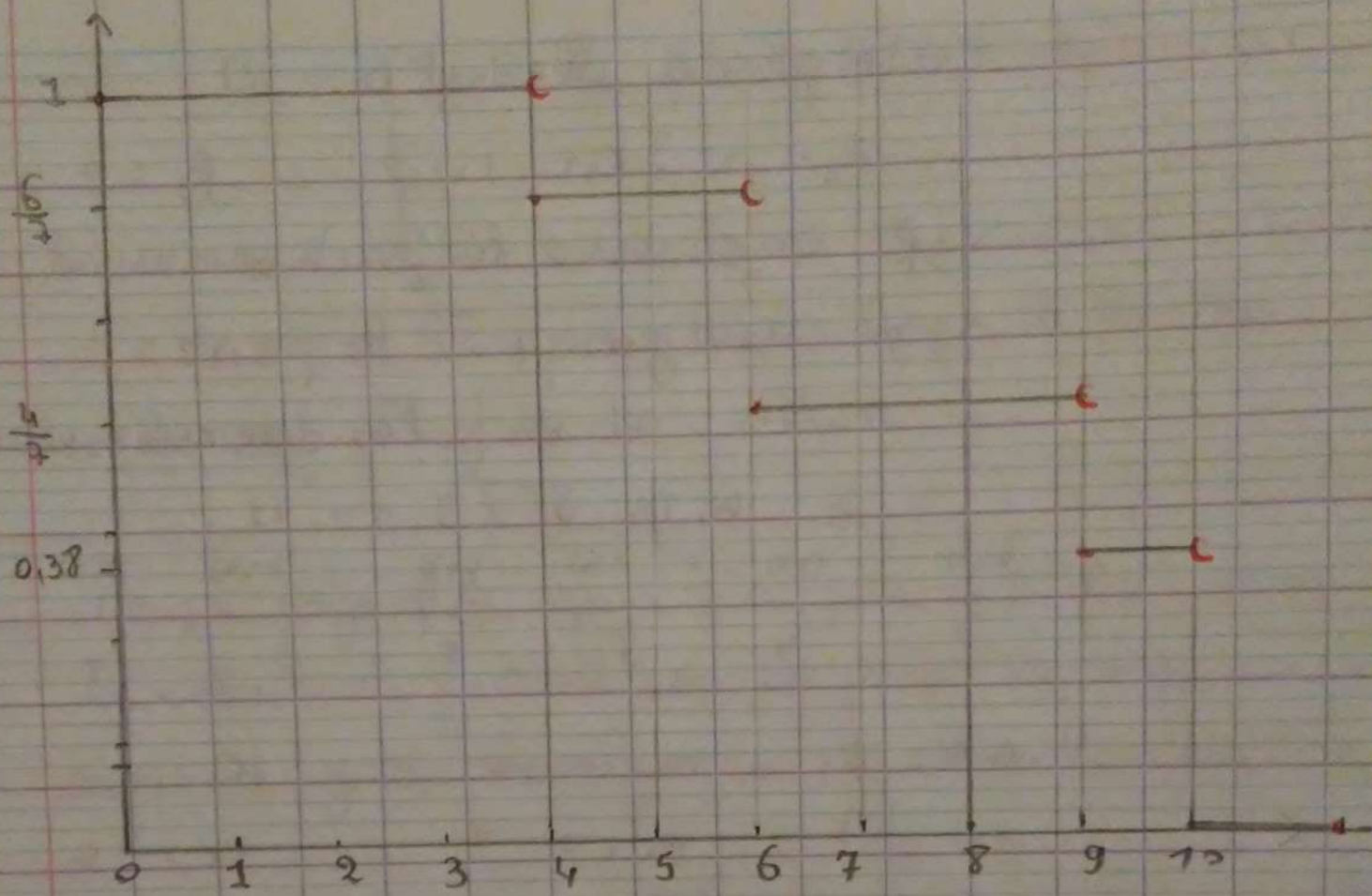
$$S_{KH}(t) = \frac{4}{7} \left(1 - \frac{m_3}{n_3}\right)^0 = \frac{4}{7} \quad t \in [8, 9[$$

$$S_{KH}(t) = \frac{4}{7} \left(1 - \frac{m_4}{n_4}\right)^1 = \frac{8}{7} \quad t \in [9, 10[$$

$$S_{KH}(t) = \frac{8}{2 \cdot 1} \left(1 - \frac{m_5}{n_5}\right) = \frac{4}{2 \cdot 1} \quad t \in [10, 13[$$

$$S_{KH}(t) = \frac{4}{2 \cdot 1} \left(1 - \frac{m_6}{n_6}\right) = \frac{4}{2 \cdot 1} \cdot 0 = 0 \quad t \in [13, 13]$$





$$S_n(t) = \frac{1}{n} \sum_{(T_i > t)} 1 = \underbrace{S_{n_1}(t)}_{\text{obs}} + \underbrace{S_{n_2}(t)}_{\text{censure}} \stackrel{?}{\rightarrow} S(t)$$