

# Analiza danych

prof. UAM dr hab. Tomasz Górecki

tomasz.gorecki@amu.edu.pl  
<http://drizzt.home.amu.edu.pl>

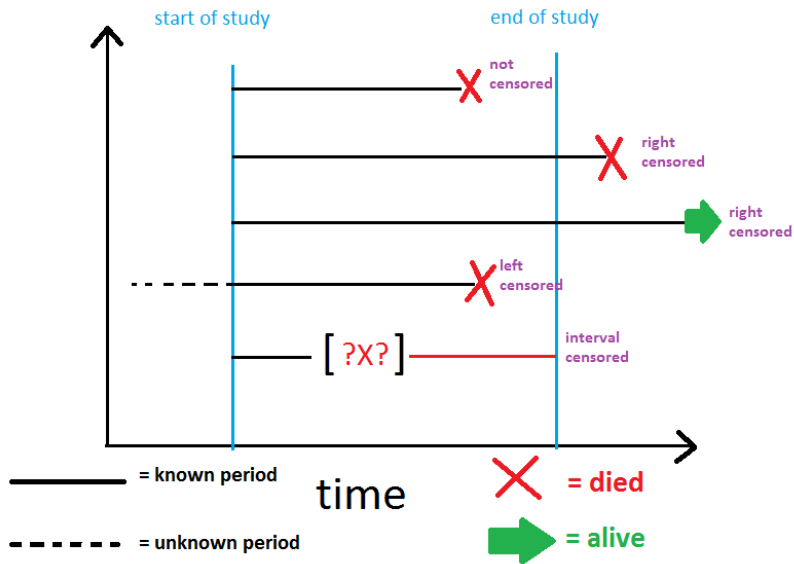
Zakład Statystyki Matematycznej i Analizy Danych  
Wydział Matematyki i Informatyki  
Uniwersytet im. Adama Mickiewicza w Poznaniu



## Wprowadzenie

Przedmiotem badania **analizy przeżycia** (*ang. survival analysis*) jest czas jaki upływa od początku obserwacji do wystąpienia określonego zdarzenia, które jednoznacznie kończy obserwację na danej jednostce eksperymentalnej. Analiza przeżycia wywodzi się wprawdzie z badań medycznych (na co wskazuje nazwa), lecz znajduje również zastosowanie w innych naukach (np. analiza niezawodności w naukach technicznych). Charakterystyczne dla analizy przeżycia są tzw. **dane cenzorowane** (*ang. censored observations*) inaczej **ucięte** (oznaczamy  $70^+$ ), o których wiadomo, że zdarzenie nie nastąpiło aż do momentu zakończenia obserwacji (np. pacjenci wypisani ze szpitala).

# Wprowadzenie



## Wprowadzenie

Głównym obiektem badawczym jest tzw. **funkcja przeżycia** (ang. *survival function*)  $S(t)$ , która określa prawdopodobieństwo, że osoba przeżyje dłużej niż pewien przyjęty czas  $t$ , czyli

$$S(t) = \mathbb{P}(T > t) = \int_t^{\infty} f(u)du = 1 - F(t),$$

gdzie  $T$  jest absolutnie ciągłą zmienną losową określającą czas życia o funkcji gęstości  $f$  oraz dystrybuancie  $F$  i przyjmującą wartości z przedziału  $[0, \infty)$

## Estymator Kaplana-Meiera

Najpopularniejszym estymatorem funkcji przeżycia jest **estymator KAPLANA-MEIERA** (ang. *Kaplan's Meier estimator*) postaci:

$$\hat{S}(t) = \prod_{t_i \leq t} \frac{n_i - d_i}{n_i},$$

gdzie  $n_i$  jest liczbą obiektów, która dożyła momentu  $t_i$ , natomiast  $d_i$  jest liczbą śmierci w momencie  $t_i$ . Ważną jego zaletą jest uwzględnianie obserwacji cenzorowanych.

## Estymator Kaplana-Meiera – przykład

Znane są czasy przeżycia (w dniach) dla 16 pacjentów, którzy przeszli zabieg usunięcia guza mózgu: 28, 49, 54, 80, 80, 102<sup>+</sup>, 120, 120<sup>+</sup>, 120<sup>+</sup>, 167, 200, 200, 200<sup>+</sup>, 340, 500, 500<sup>+</sup>.

Znajdziemy dla tych danych estymator KM funkcji przeżycia. Założmy wpraw, że nie uwzględniamy cenzorowania danych.

$i$	$t_i$	$n_i$	$d_i$	$\frac{n_i - d_i}{n_i}$	$\hat{S}(t = t_i)$
1	28	16	1	0.9375	0.9375
2	49	15	1	0.9333	0.8750
3	54	14	1	0.9286	0.8125
4	80	13	2	0.8462	0.6875
5	102	11	1	0.9091	0.6250
6	120	10	3	0.7000	0.4375
7	167	7	1	0.8571	0.3750
8	200	6	3	0.5000	0.1875
9	340	3	1	0.6667	0.1250
10	500	2	2	0.0000	0.0000

## Estymator Kaplana-Meiera – przykład

Z kolei jeśli uwzględnimy, że część obserwacji była cenzorowana:

$i$	$t_i$	$n_i$	$d_i$	$\frac{n_i - d_i}{n_i}$	$\hat{S}(t = t_i)$
1	28	16	1	0.9375	0.9375
2	49	15	1	0.9333	0.8750
3	54	14	1	0.9286	0.8125
4	80	13	2	0.8462	0.6875
5	120	10	1	0.9000	0.6188
6	167	7	1	0.8571	0.5304
7	200	6	2	0.6667	0.3536
8	340	3	1	0.6667	0.2357
9	500	2	1	0.5000	0.1179

Dane: 28, 49, 54, 80, 80, 102<sup>+</sup>, 120, 120<sup>+</sup>, 120<sup>+</sup>, 167, 200, 200, 200<sup>+</sup>, 340, 500, 500<sup>+</sup>.

## Model Coxa

Estymator KM pozwala na graficzną prezentację krzywej przeżycia oraz porównanie takich krzywych dla kilku populacji. Nie jest jednak możliwe za jego pomocą opisanie zależności funkcji przeżycia od innych zmiennych objaśniających. Wydawać by się mogło, że do tego zagadnienia powinniśmy zastosować poznane już modele regresji wielokrotnej. Nie jest to jednak możliwe z dwóch powodów: czas przeżycia przeważnie nie ma rozkładu normalnego oraz występuje problem wykorzystania obserwacji cenzorowanych.



## Model Coxa

Najpopularniejszym modelem regresji wykorzystywanym w tym przypadku jest **model proporcjonalnego hazardu COXA** (*ang. Cox proportional hazard model*). Zdefiniujmy wpieryw **funkcję hazardu** (*ang. hazard function*):

$$\lambda(t) = \lim_{dt \rightarrow 0} \frac{\mathbb{P}(t \leq T \leq t + dt)}{dt \cdot S(t)} = \frac{f(t)}{S(t)} = \frac{S'(t)}{S(t)}.$$

Funkcja hazardu nie jest prawdopodobieństwem, ale miarą ryzyka śmierci w chwili  $t$  (im większy hazard tym większe ryzyko).

Każda nieujemna funkcja  $\lambda$ , dla której zachodzi

$$\int_0^{\infty} \lambda(t) dt = \infty$$

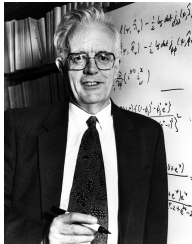
może być funkcją hazardu.

## Model Coxa

Model proporcjonalnego hazardu COXA ma postać:

$$\lambda(t | \mathbf{X}) = \lambda_0(t)e^{\mathbf{X}\beta},$$

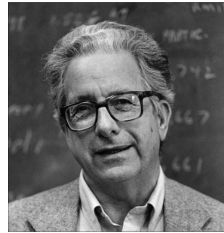
gdzie  $\lambda_0(t)$  jest **hazardem zerowym** (ang. *baseline hazard function*) – hazard, gdy wszystkie zmienne niezależne są równe zero. W modelu tym nie zakładamy nic o postaci funkcji hazardu (takie założenie można wprowadzić uzyskując parametryczne modele proporcjonalnego hazardu). Jeśli założymy, że elementy wektora parametrów nie zależą od czasu, to mamy do czynienia z modelem proporcjonalnego hazardu COXA.



Sir David Roxbee Cox  
(1924-2022)



Edward Lynn Kaplan  
(1920-2006)



Paul Meier  
(1924-2011)

## Bibliografia



Cox, D.R. (1972). Regression models and life tables. *Journal of the Royal Statistical Society: Series B* 34(2):187-220.



Kaplan, E.L., Meier, P. (1958). Nonparametric estimation from incomplete observations. *Journal of the American Statistical Association* 53(282):457-481.

