

# **A flexible approach to time-to-event data analysis using case-base sampling**

---

Jesse Islam

July 11, 2019

## Motivating example

- Meet Justin.

## Motivating example

- Meet Justin.
  - Age: 56

## Motivating example

- Meet Justin.
  - Age: 56
  - Worried about his prostate.

# Motivating example

- Meet Justin.
  - Age: 56
  - Worried about his prostate.
  - What is Justin's two year risk of death due to prostate cancer?

## Popular methods in time-to-event analysis

- In disease etiology, we tend to make use of the proportional hazards hypothesis .

## Popular methods in time-to-event analysis

- In disease etiology, we tend to make use of the proportional hazards hypothesis .
  - Cox Regression

# Popular methods in time-to-event analysis

- In disease etiology, we tend to make use of the proportional hazards hypothesis .
  - Cox Regression
- When we want the absolute risk:



# Popular methods in time-to-event analysis

- In disease etiology, we tend to make use of the proportional hazards hypothesis .
  - Cox Regression
- When we want the absolute risk:
  - Parametric models

# Popular methods in time-to-event analysis

- In disease etiology, we tend to make use of the proportional hazards hypothesis .
  - Cox Regression
- When we want the absolute risk:
  - Parametric models
  - Breslow estimator

# Motivations for a new method

- Julien and Hanley found that survival analysis rarely produces prognostic functions, even though the software is widely available in cox regression packages. [1]

# Motivations for a new method

- Julien and Hanley found that survival analysis rarely produces prognostic functions, even though the software is widely available in cox regression packages. [1]
- They believe the stepwise nature is the reason, as it reduces interpretability. [1]

# Motivations for a new method

- Julien and Hanley found that survival analysis rarely produces prognostic functions, even though the software is widely available in cox regression packages. [1]
- They believe the stepwise nature is the reason, as it reduces interpretability. [1]
- Want to easily model non-proportional hazards. [1]

# Motivations for a new method

- Julien and Hanley found that survival analysis rarely produces prognostic functions, even though the software is widely available in cox regression packages. [1]
- They believe the stepwise nature is the reason, as it reduces interpretability. [1]
- Want to easily model non-proportional hazards. [1]
- A streamlined approach for reaching a **smooth absolute risk** curve. [1]

## Dr. Cox's perspective

**Reid:** How do you feel about the cottage industry that's grown up around it [the Cox model]?

**Cox:** Don't know, really. In the light of some of the further results one knows since, I think I would normally want to tackle problems parametrically, so I would take the underlying hazard to be a Weibull or something. I'm not keen on nonparametric formulations usually.

**Reid:** So if you had a set of censored survival data today, you might rather fit a parametric model, even though there was a feeling among the medical statisticians that that wasn't quite right.

**Cox:** That's right, but since then various people have shown that the answers are very insensitive to the parametric formulation of the underlying distribution [see, e.g., Cox and Oakes, Analysis of Survival Data, Chapter 8.5]. And if you want to do things like predict the outcome for a particular patient, it's much more convenient to do that parametrically.

# European Randomized Study of Prostate Cancer Screening (ERSPC) Data

- ~150 000 men ages 55-69.

## The European Randomized Study of Screening for Prostate Cancer – Prostate Cancer Mortality at 13 Years of Follow-up

Fritz H. Schröder<sup>1</sup>, Jonas Hugosson<sup>2</sup>, Monique J. Roobol<sup>1</sup>, Teuvo L.J. Tammela<sup>3</sup>, Marco Zappa<sup>4</sup>, Vera Nelen<sup>5</sup>, Maciej Kwiatkowski<sup>6,7</sup>, Marcos Lujan<sup>8,9</sup>, Lissa Määttänen<sup>10</sup>, Hans Lilja<sup>11,12,13</sup>, Louis J. Denis<sup>14</sup>, Franz Recker<sup>6</sup>, Alvaro Paez<sup>15,16</sup>, Chris H. Bangma<sup>1</sup>, Sigrid Carlsson<sup>2,11</sup>, Donella Puliti<sup>4</sup>, Arnauld Villers<sup>17</sup>, Xavier Rebillard<sup>18</sup>, Matti Hakama<sup>10,19</sup>, Ulf-Hakan Stenman<sup>20</sup>, Paula Kujala<sup>21</sup>, Kimmo Taari<sup>22</sup>, Gunnar Aus<sup>23</sup>, Andreas Huber<sup>24</sup>, Theo van der Kwast<sup>25</sup>, Ron H.N. van Schaik<sup>26</sup>, Harry J. de Koning<sup>27</sup>, Sue M. Moss<sup>28</sup>, Anssi Auvinen<sup>19</sup>, and for the ERSPC Investigators



# European Randomized Study of Prostate Cancer Screening (ERSPC) Data

- ~150 000 men ages 55-69.
- Examined effects screening has on death due to prostate cancer.

## The European Randomized Study of Screening for Prostate Cancer – Prostate Cancer Mortality at 13 Years of Follow-up

Fritz H. Schröder<sup>1</sup>, Jonas Hugosson<sup>2</sup>, Monique J. Roobol<sup>1</sup>, Teuvo L.J. Tammela<sup>3</sup>, Marco Zappa<sup>4</sup>, Vera Nelen<sup>5</sup>, Maciej Kwiatkowski<sup>6,7</sup>, Marcos Lujan<sup>8,9</sup>, Lissa Määttänen<sup>10</sup>, Hans Lilja<sup>11,12,13</sup>, Louis J. Denis<sup>14</sup>, Franz Recker<sup>6</sup>, Alvaro Paez<sup>15,16</sup>, Chris H. Bangma<sup>1</sup>, Sigrid Carlsson<sup>2,11</sup>, Donella Puliti<sup>4</sup>, Arnaud Villers<sup>17</sup>, Xavier Rebillard<sup>18</sup>, Matti Hakama<sup>10,19</sup>, Ulf-Hakan Stenman<sup>20</sup>, Paula Kujala<sup>21</sup>, Kimmo Taari<sup>22</sup>, Gunnar Aus<sup>23</sup>, Andreas Huber<sup>24</sup>, Theo van der Kwast<sup>25</sup>, Ron H.N. van Schaik<sup>26</sup>, Harry J. de Koning<sup>27</sup>, Sue M. Moss<sup>28</sup>, Anssi Auvinen<sup>19</sup>, and for the ERSPC Investigators

## ERSPC Data

```
head(casebase::ERSPC)
```

PatientID	ScrArm	Follow.Up.Time	DeadOfPrCa
1	1	0.003	0
2	0	1.038	1
3	1	7.966	1
4	0	11.975	1
5	1	14.910	0

- Using the ERSPC dataset and casebase, we will determine Justin's absolute risk for death by prostate cancer.

1. Clever sampling.

# Casebase Overview

1. Clever sampling.
2. Implicitly deals with censoring.

# Casebase Overview

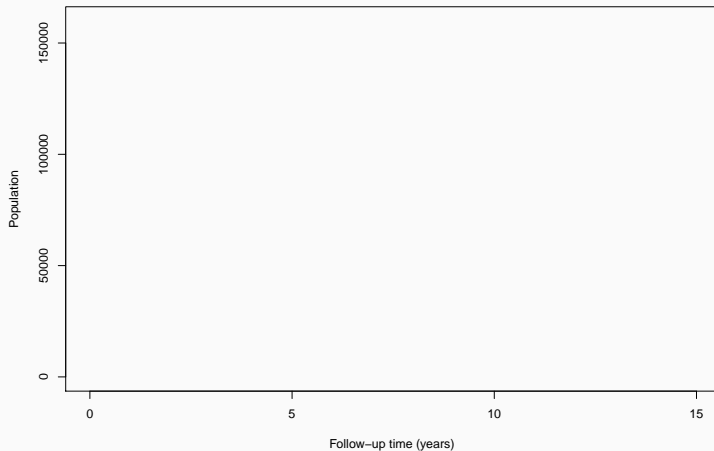
1. Clever sampling.
2. Implicitly deals with censoring.
3. Allows a parametric fit using *logistic regression*.

1. Clever sampling.
  2. Implicitly deals with censoring.
  3. Allows a parametric fit using *logistic regression*.
- Casebase is parametric, and allows different parametric fits by incorporation of the time component.

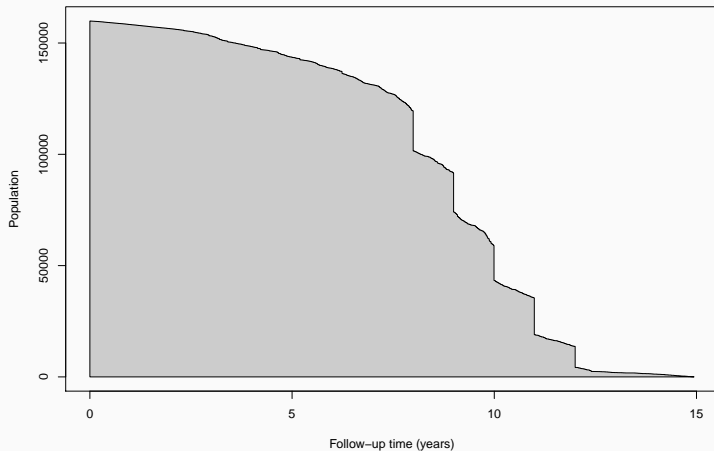
1. Clever sampling.
  2. Implicitly deals with censoring.
  3. Allows a parametric fit using *logistic regression*.
- Casebase is parametric, and allows different parametric fits by incorporation of the time component.
  - Package contains an implementation for generating *population-time* plots.



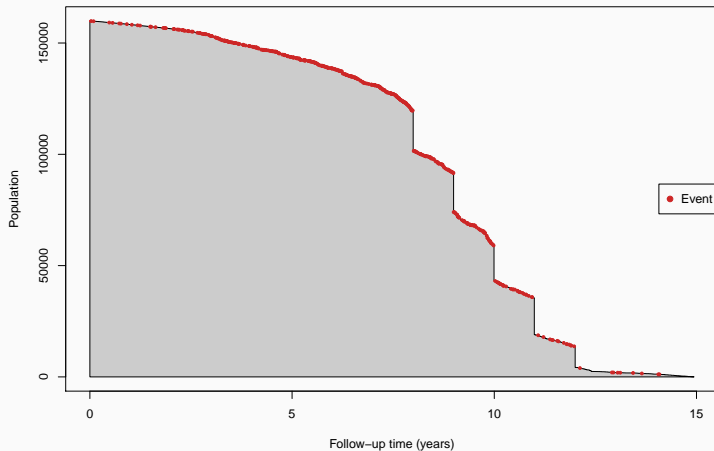
# Casebase: Sampling



# Casebase: Sampling

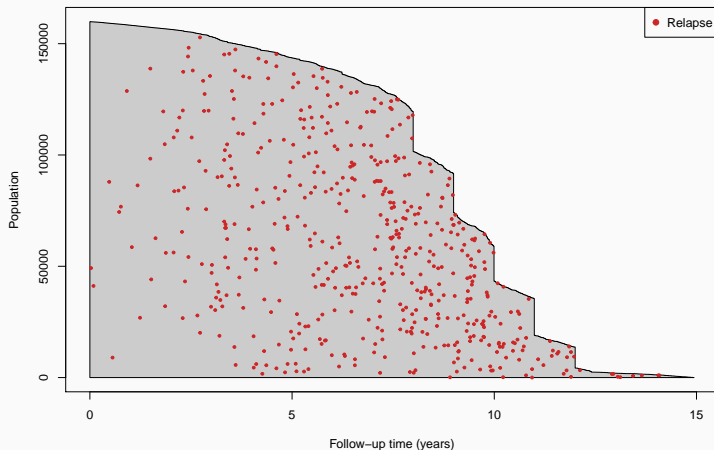


# Casebase: Sampling

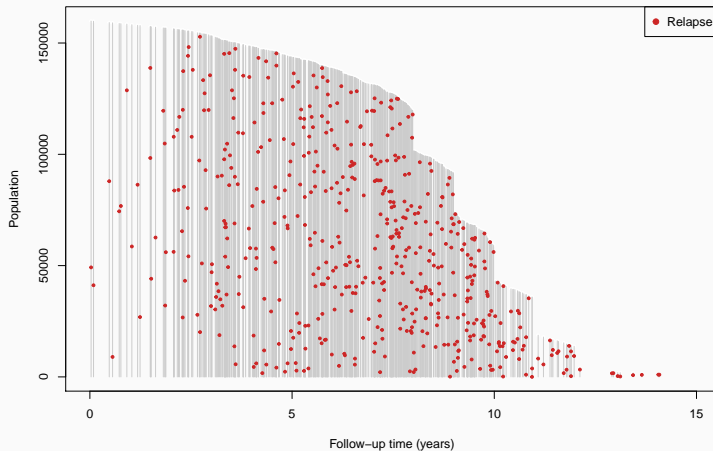


# Casebase: Sampling

```
casebase::popTime(Data,Event,Time)
```



# Casebase: Sampling



- We can now fit models of the form:

$$\log(h(t; \alpha, \beta)) = g(t; \alpha) + \beta X$$

- We can now fit models of the form:

$$\log(h(t; \alpha, \beta)) = g(t; \alpha) + \beta X$$

- By changing the function  $g(t; \alpha)$ , we can model different parametric families easily:

## Casebase: Parametric models

*Exponential:*  $g(t; \alpha)$  is equal to a constant

```
casebase::fitSmoothHazard(status ~ X1 + X2)
```

*Gompertz:*  $g(t; \alpha) = \alpha t$

```
casebase::fitSmoothHazard(status ~ time + X1 + X2)
```

*Weibull:*  $g(t; \alpha) = \alpha \log(t)$

```
casebase::fitSmoothHazard(status ~ log(time) + X1 + X2)
```



# Death by prostate cancer: hazard ratios

```
casebase::fitSmoothHazard(DeadOfPrCa~ log(Follow.Up.Time)+  
                           ScrArm, data=ERSPC, ratio = 100)
```

Call:

```
glm(formula = formula, family = binomial, data = sampleData)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-0.2693	-0.1715	-0.1348	-0.0908	4.5189

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-9.46535	0.15812	-59.862	<2e-16 ***
log(Follow.Up.Time)	1.08124	0.08264	13.084	<2e-16 ***
ScrArm	-0.20833	0.08859	-2.352	0.0187 *

---

signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 6059.0 on 54539 degrees of freedom  
Residual deviance: 5794.1 on 54537 degrees of freedom  
AIC: 5800.1

Number of Fisher Scoring iterations: 8

## ERSPC Hazard comparison

Model	Hazard Ratio	Std.Error
Cox	0.801	1.092
Gompertz	0.802	1.093
Exponential	0.810	1.092
Weibull	0.797	1.093

- We have a bunch of different parametric hazard models now.

- We have a bunch of different parametric hazard models now.
- To get the absolute risk, we need to evaluate the following equation in relation to the hazard:

$$CI(x, t) = 1 - e^{-\int_0^t h(x, u) du}$$

- We have a bunch of different parametric hazard models now.
- To get the absolute risk, we need to evaluate the following equation in relation to the hazard:

$$CI(x, t) = 1 - e^{-\int_0^t h(x, u) du}$$

- $CI(x, t)$  = Cumulative Incidence (Absolute Risk)

# Absolute Risk

- We have a bunch of different parametric hazard models now.
- To get the absolute risk, we need to evaluate the following equation in relation to the hazard:

$$CI(x, t) = 1 - e^{-\int_0^t h(x, u) du}$$

- $CI(x, t)$  = Cumulative Incidence (Absolute Risk)
- $h(x, u)$  = Hazard function

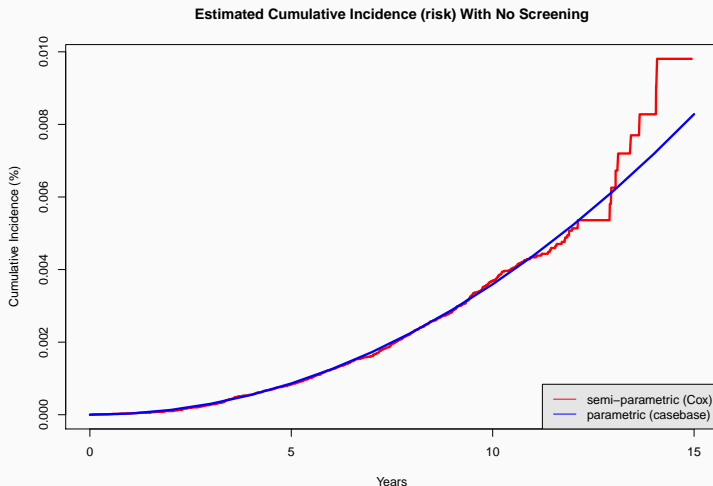
- We have a bunch of different parametric hazard models now.
- To get the absolute risk, we need to evaluate the following equation in relation to the hazard:

$$CI(x, t) = 1 - e^{-\int_0^t h(x, u) du}$$

- $CI(x, t)$  = Cumulative Incidence (Absolute Risk)
- $h(x, u)$  = Hazard function
- Lets use the weibull hazard

# Casebase: Absolute Risk comparison

```
casebase::absoluteRisk(fit, time=5, covariate_profile)
```





- Casebase sampling implicitly incorporates censoring and permits the use of GLMs and the tools associated with them

- Casebase sampling implicitly incorporates censoring and permits the use of GLMs and the tools associated with them
- The casebase package contains tools to generate:

- Casebase sampling implicitly incorporates censoring and permits the use of GLMs and the tools associated with them
- The casebase package contains tools to generate:
  - Population-Time plots

- Casebase sampling implicitly incorporates censoring and permits the use of GLMs and the tools associated with them
- The casebase package contains tools to generate:
  - Population-Time plots
  - Hazard functions

- Casebase sampling implicitly incorporates censoring and permits the use of GLMs and the tools associated with them
- The casebase package contains tools to generate:
  - Population-Time plots
  - Hazard functions
  - Absolute Risk

- Casebase sampling implicitly incorporates censoring and permits the use of GLMs and the tools associated with them
- The casebase package contains tools to generate:
  - Population-Time plots
  - Hazard functions
  - Absolute Risk
  - Casebase can deal with competing risks.

## References 1

1. Hanley, James A, and Olli S Miettinen. 2009. "Fitting Smooth-in-Time Prognostic Risk Functions via Logistic Regression." *The International Journal of Biostatistics* 5 (1).
2. Olli presentation slides?
3. Saarela, Olli. 2015. "A Case-Base Sampling Method for Estimating Recurrent Event Intensities." *Lifetime Data Analysis*. Springer, 1–17

## References 2

- 4.Schroder FH, et al., for the ERSPC Investigators.Screening and Prostate-Cancer Mortality in a Randomized European Study. *N Engl J Med* 2009;360:1320-8.
- 5.Scrucca L, Santucci A, Aversa F. Competing risk analysis using R: an easy guide for clinicians. *Bone Marrow Transplant.* 2007 Aug;40(4):381-7. doi: 10.1038/sj.bmt.1705727.
- 6.Turgeon, M. (2017, June 10). Retrieved May 05, 2019, from <https://www.maxturgeon.ca/slides/MTurgeon-2017-Student-Conference.pdf>



Tutorial:

<http://sahirbhatnagar.com/casebase/>

Slides:

<https://github.com/Jesse-Islam/UseR-CaseBase-Presentation>

Questions?



- Current methods:

- Current methods:
  - Fine-Gray

- Current methods:
  - Fine-Gray
  - Kaplan-Meier

- Current methods:
  - Fine-Gray
  - Kaplan-Meier
- Proposed method:

- Current methods:
  - Fine-Gray
  - Kaplan-Meier
- Proposed method:
  - Case-Base

## Competing Risks: Data

- Two diseases:

```
head(casebase::bmtcrr)
```

D	Status	ftime
ALL	2	0.67
AML	1	9.50
ALL	0	131.77
ALL	2	24.03



## Competing Risks: Data

- Two diseases:
  - Acute Lymphoblastic leukemia (ALL)

```
head(casebase::bmtcrr)
```

D	Status	ftime
ALL	2	0.67
AML	1	9.50
ALL	0	131.77
ALL	2	24.03

## Competing Risks: Data

- Two diseases:
  - Acute Lymphoblastic leukemia (ALL)
  - Acute Myeloblastic leukemia (AML)

```
head(casebase::bmtcrr)
```

D	Status	ftime
ALL	2	0.67
AML	1	9.50
ALL	0	131.77
ALL	2	24.03

## Competing Risks: Data

- Two diseases:
  - Acute Lymphoblastic leukemia (ALL)
  - Acute Myeloblastic leukemia (AML)
- Contains a competing event.

```
head(casebase::bmtcrr)
```

D	Status	ftime
ALL	2	0.67
AML	1	9.50
ALL	0	131.77
ALL	2	24.03

## Competing Risks: Absolute Risk

```
fit_cb <- casebase::fitSmoothHazard(Status ~ ftime  
                                     + ... , data =  
                                     bmtcrr)  
risk_cb <- absoluteRisk(fit_cb, Time, Newdata)
```

# Competing Risks: Absolute Risk

