



Statistical dilemmas in personalized risk prediction: the Estonian Biobank example

Krista Fischer

Institute of Mathematics and Statistics, Institute of Genomics, University of Tartu

Estonian Academy of Sciences





Outline

- What is risk?
- Risk of death
- Basic methods for survival analysis
- Risk of a disease and competing risks
- Estonian Biobank and solutions for personalized medicine: an overview
- Risk prediction for common complex diseases in the Estonian Biobank



What is probability?

A quiz (D. Spiegelhalter)

https://tinyurl.com/yyeavsa4



How good we are at estimating probabilities?

How was the score calculated?

5 0 0 6 9 -11 7 16 -24 8 21 -39		Confidence	Score if you are right	Score if you are wrong
6 9 -11 7 16 -24 8 21 -39	The Prior score	5	0	0
7 16 -24 8 21 -39 0 0 56	The brief score.	6	9	-11
8 21 -39		7	16	-24
		8	21	-39
9 24 -56		9	24	-56
10 25 -75	M/hy^2	10	25	-75

- Know all the answers and are 100% confident in them (score 10): **25** points per question.
- Choose the answer randomly, but still choose the confidence score 10: on average, you will get half of the answers correct and half will be wrong: (25-75)/2=-25
 - Know the answers with the probability of 80% and choosing 8 as the confidence.
 Average: 0,8*21-0,2*39= 9 points.
 - Know the answers with the probability of 50% and choosing 8 as the confidence.
 Average: 0,5*21-0,5*39= -9 points, etc.

Too much confidence will lead to a negative score: very often people overestimate their probability of knowing the answer!

nature

Explore content v About the journal v Publish with us v Subscribe

nature > essay > article

ESSAY 16 December 2024 Correction <u>18 December 2024</u>

Why probability probably doesn't exist (but it is useful to act like it does)

All of statistics and much of science depends on probability – an astonishing achievement, considering no one's really sure what it is.

By David Spiegelhalter

Although a bit confusing, it is still useful to use probabilities to speak about risks (or do we have a better alternative?)

- What is the risk of death?
- The probability of death (at some point) is 100% for each and every one of us.
- We always have to keep a time axis in mind death in next 10 years, before the age of 80, before 2050, ...

Although a bit confusing, it is still useful to use probabilities to speak about risks (or do we have a better alternative?)

- What is the risk of death?
- The probability of death (at some point) is 100% for each and every one of us.
- We always have to keep a time axis in mind death in next 10 years, before the age of 80, before 2050, ...

A bit of history...

Wilhelm Lexis (1837-1914)

- German statistician, economist, and social scientist
- Pioneer of the analysis of demographic time series
- Professor at the University of Dorpat (now: Tartu University) in 1874-1876 (chair of geography, ethnography and statistics)



Introduction to the Theory of Population Statistics by W. Lexis

- One of his most famous works
- Published 150 years ago
- Still known and cited by demographers, epidemiologists and also actuarians

	1878.2555.
	EINLEITUNG
	IN DIE
	THEORIE
	DER
BEVÖ	LKERUNGSSTATISTIK
	VON
DR. D	W. LEXIS ER STAATSWISSENSCHAFTEN UND DER PHILOSOPHIE, 0. PROFESSOR DER STATISTIK IN DORPAT.
	STRASSBURG
	КАПЬ J. ТПÜВИЕП 1875.

The Lexis diagram (1875)

- Horizontal axis: birth cohorts (year of birth)
- Vertical axis: age
- Diagonal lines: calendar years



Lexis, 1875, Figure 1

A modern version of the Lexis diagram

- Introduced by R. Pressat (1961), now used widely in demography and epidemiology
- Basis of age-period-cohort modeling
- The risk can vary across each of the three time scales
- The problem if identifiablity: linear dependence

```
Age = Calendar time – Birth time
```



Mai-Britt Meriloo, Bachelor's thesis 2025

A conference advert

150 years of Lexis diagram

- Tartu, October 17-18, 2025
- Organized by the Estonian Statistical Association
- (17.10: presentations in English, 18.10: Estonian)

The risk of death: time-to-event analysis approach

T: time to (age at) death.

```
Survival function: P(T > t)
```

- Probability of survival up to time (age) t.
- Easy to estimate (classical probability), if we have complete data: a cohort of individuals followed from birth to death.
- Usually, we don't have such cohorts!

The idea of Kaplan-Meier estimator (and other similar estimators)

A fictional study: duration 5 years. Suppose we start with individuals who are either 65, 70 or 75 in the beginning and record 5-year mortality.

Age	Ν	Alive in 5 years	5-year mortality	5-year survival
65	1000	850	0.15	0.85
70	700	560	0.2	0.8
75	600	450	0.25	0.75

We can estimate the conditional probabilities P(T > 70|T > 65) (5-year survival among 65-year-olds), P(T > 75|T > 70) and P(T > 80|T > 75).

Assuming that the survival rates stay the same in time, we can estimate the probability of a 65-year old to live up to 80, by combining these together:

 $P(T > 80|T > 65) = P(T > 80|T > 75) \cdot P(T > 75|T > 70) \cdot P(T > 70|T > 65).$

Thus, the probability here is $P(T > 80 | T > 65) = 0.85 \cdot 0.8 \cdot 0.75 = 0.51$.

Survival curves based on one year of data (Statistics Estonia)



Useful for population statistics, but useless for individual survival prediction!

(Think of the Lexis diagram)



Calendar time

The Estonian Biobank: from population-based biobank to personalized medicine

Human Resear	Genes rch Act	ι	→ Jniversity	, ,	Vision of ersonaliz medicine	f ed	Retur of results	F		
Idea	EstBB Iaunch		of Tartu			Pilot projects		100,000 Genome Project		
1999 20	00 2001	2004	2007	2010	2013	2015	2017	2018	2019	2020
E	Biobank participa	ants 10 OC	00	50 000						200 000
- A	2						Nationa	al Personalizec	l Medicine	Project

Prof. Andres Metspalu



211 000+ biobank participants

Health records, diet, physical activity etc

DNA, plasma and cell samples



Estonian Biobank: number of participants recruited per year



Survival curves based on cohort (biobank) data



Survival in the Estonian Biobank cohort

This reflects average survival across different recruitment times (2002-2019)

Survival curves based on cohort (biobank) data



Age

Survival in the two cohorts of the Estonian Biobank

Very much dependent on the cohort!

What data could be used for individual risk prediction and how?

Example: NMR-biomarkers and mortality (work with Mara Delesa-Velina)



Difficult to use in the individual risk prediction

The risk score developed in the "old" cohort, predictions in the new cohort (validation set)

An alternative idea

Survival-based biological

age: age, where the average survival probability in the population equals to the individual's current survival probability, given his/her covariate profile.

 How to compute? Parametric survival modeling (Gompertz, Weibull...)

Survival-based biological age



Biological age estimates



(Mara Delesa-Velina)

Effects of risk factors



The risk of a disease

Are we interested in...

- Probability of ever getting the disease?
- Probability of getting the disease in X (5, 10, ...) years?
- Probability of getting the disease before age A (60, 70, ...)?

Probability of disease is often expressed as $P(T \le t)$ and estimated as cumulative incidence:

Risk prediction for common complex diseases (the biobank view)



Personalized risk prediction for common complex diseases has existed long before biobanks

Common risk prediction algorithms:

- Coronary heart disease: SCORE, PCE, QRISK
- Diabetes: FinDRisc, QDiabetes, ...
- Cancer: QCANCER...
- Our aim is to add the genetic component to the established risk algorithms (calibrated for the Estonian data)

Is it worth adding the genetic component to the risk prediction algorithm?

- Genetic risk component summarized as the Polygenic Risk Score (PRS)
- Need to select the best PRS among alternatives (PGS Catalogue!)
- Does it explain a meaningful amount of variability?
- Need to validate in the (sub)cohort that was not included in the PRS development process!

MetaGRS (combining 2 PRS-s) for incident **Breast Cancer** in the Estonian Biobank cohort



Cumulative incidence: $P(T \le t)$

Läll et al, 2019, BMC Cancer

Type 2 Diabetes

Also:

Läll 2017,

Gen.Med



PRS group - <40% - 40-60% - 60-80% - 80-95% - >95%



Coronary Artery Disease



Again, we see different risks in the two cohorts

(Tuuli Puusepp)

PRS percentile — Bottom 10% — 10%–90% — Top 10%

https://www.medrxiv.org/content/10.1101/2025.04.02.25324383v1

Challenges

- Need to move from effect estimation and testing to absolute risk prediction
- What is the best way to communicate personalized risks?

Main steps of developing an algorithm for risk prediction and communication: the example of cardiovascular risk prediction in the Estonian Biobank

- Statistical modeling to assess the effects of risk factors and to develop the optimal predictive model
- Developing a predictive tool for absolute risk
- Finding the best way to communicate the risks

Assessing the effect of the PRS in the Estonian Biobank (33082 men and 74629 women of age 18-80) ...



Kaplan-Meier curves illustrate the average effect size, but they are...

- not adjusted for other risk factors
- not easily implemented for out-of-sample predictions (nonparametric)

PRS: metaGRS from Innouye et al. (JACC, 2018)

The proportional hazards model

The hazard function:

 $h(t) = \lim_{dt \to 0} P(t \le T < t + dt | T > t) / dt$



... is a multiplicative model for hazard: when a covariate changes by a constant, the hazard is multiplied by a constant

The Cox proportional hazards model

$$h(t|X_1,\ldots,X_k) = h_0(t)e^{X_1\beta_1+\cdots+X_k\beta_k}$$

 A semiparametric model – does not use the real survival times, but the ranks of survival time and information on the individuals at risk at each event time.

 The baseline hazard h₀ will not be estimated – thus the model does only provide estimates of the covariate effects, but does not allow direct prediction of hazards or survival times



```
Sir David Cox
(Oxford, UK)
Paper on prop.haz
models in 1972
```

Estimated effects of risk factors on incident CHD (Hazard Ratios with 95%CI)



We see the adjusted estimates and their range of uncertainty (CI-s), but ...

- HR-s are not so easily understood by general public
- need baseline hazard estimates for actual risk prediction

Cumulate incidence estimates from the Cox model...



... good to illustrate the effects of risk factors on absolute risks, but...

- it still relies on nonparametric baseline hazard estimation
- does not account for the fact that an individual does not have CVD at baseline (varying age)

age



Figure 4. Reclassification of individuals initially categorized as intermediate risk for 5year CVD incidence using the conventional model. Arrows indicate the movement of individuals between categories, with corresponding percentages representing the proportion of individuals reclassified.

Tuuli Puusepp et al. (2025)

What do we actually need to estimate for individual risk prediction?

- Feedback on risks is relevant for the individuals who are currently disease-free
- Often, a 10-year risk is a meaningful quantity to be estimated for risk stratification/feedback purposes
- Instead of the popular Cox model, parametric modeling of time-toevent censored outcomes deserves more attention, providing more straightforward tools for risk prediction

An alternative: a Weibull model for 10-year risk

Predicted 10-year risks for men in the Estonian Biobank



Current age

10-year risk:

 $P(T_i < 10) = 1 - e^{-(10\lambda_i)^{\gamma}}$

From the model

Geneetiline risk võrreldes rahvastikuga

Implementation in the portal





~

tasakaalustatud

kehakaal

Koguriski muutumine aastatega





Salvestamiseks mine profiilile \rightarrow

Further challenges: work in progress The current model is not ideal for everyone...

Predicted 10-year risks for men in the Estonian Biobank



Challenges:

- Feedback to old people (competing risks)
- Feedback to young people (very low risks)



Accounting for competing risks (men, EstBB)



- Ignoring competing risks biases the risk predictions in old age
- However, communication of competing risks could be challenging
- Also not so easy to implement in a Weibull model

Incident Type 2 Diabetes (T2D) and polygenic (genetic) risk score (GRS=PRS): cumulative incidence on age scale (men with BMI=25..30) Treating death as a competing event



By age 75, in the high genetic risk group: 25% have died without a T2D diagnosis 34% have had a T2D diagnosis 41% would be alive and free of T2D (note: sum=100)

In the low genetic risk group: 32% died without T2D 12% have had a T2D diagnosis 56% alive, free of T2D

Is mortality higher in the low-PRS group?

How to do it correctly?

An example with T2D: scenarios by the age of 75 Low PRS High PRS



Application of the biological age for cardiovascular risk: heart age

Predicted 10-year risks for men in the Estonian Biobank



Real age vs heart age (men, EstBB)



Real age



Real age vs heart age (women, EstBB)

Real age

Summary

- It is not easy to speak about risks, as risk itself is confusing concept
- Risk of death is "easier" to, as death occurs only once. Still, there is the problem of reference cohort and non-identifiable age-period-cohort effects (Lexis!)
- Biological age may be easier to communicate
- Disease risk estimation requires translation from model parameters to absolute risks. For short-term prediction, in case of no considerable competing risks, there are several approaches available.
- In case of competing risks they can be taken into account in the analysis, but the communication task is still tricky.
- There are still many challenges ahead