

# A practical approach to fitting cancer survival models when data can't move across borders

Paul C Lambert<sup>1,2</sup>, Mark J Rutherford<sup>3</sup>, Tor Åge Myklebust<sup>1,4</sup>

<sup>1</sup>Cancer Registry of Norway, FHI, Norway

<sup>2</sup>Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm, Sweden

<sup>3</sup>Biostatistics Research Group, Population Health Sciences, University of Leicester, UK

<sup>4</sup>Dept. of Research and Innovation, Møre and Romsdal Hospital Trust, Ålesund, Norway

**ANCR symposium 2024, Bodø, Norway, 30th August 2024**

Slides: [pclambert.net/pdf/XXX](https://pclambert.net/pdf/XXX)

Example: [pclambert.net/software/standsurv/XXX](https://pclambert.net/software/standsurv/XXX)

# Introduction

- It is harder to share data between countries, making international comparisons more difficult.
- Here, I focus on **survival analysis**.
  - Generally need individual level data
  - Sometimes we need/want to use statistical modelling approaches (e.g. recent NORDCAN Survival Studies).
- NORDCAN.R showed how a federated approach could be applied.
  - Data analysed separately in each country
  - Aggregated/summary data sent to IARC
- Here I will explore something similar for a modelling approach.

# Joint or separate models?

- We have choices.
  - ① Fit separate models for each country.
  - ② Fit a single joint model to all countries.
- A single model can be more efficient as we can 'borrow strength' between countries.
  - However, it requires data to be in one place or to use a full federated learning approach.
- If we have large data then we are happier to fit separate models.
- A joint model with interactions between country and all covariates (and time) is equivalent to separate models.

# Options if we want a modelling approach

- ① Fit model separately in each country
  - Extract statistics of interest (e.g. 5 year relative survival)
  - Send to hub

# Options if we want a modelling approach

- ① Fit model separately in each country
  - Extract statistics of interest (e.g. 5 year relative survival)
  - Send to hub
- ② Fit model separately in each country
  - Save model object
  - Send to hub

# Options if we want a modelling approach

- ① Fit model separately in each country
  - Extract statistics of interest (e.g. 5 year relative survival)
  - Send to hub
- ② Fit model separately in each country
  - Save model object
  - Send to hub
- ③ Federated learning
  - Hub defines model
  - Parameters sent to each node
  - Aggregated model information sent back
  - Parameters updated
  - Repeat until convergence

# Options if we want a modelling approach

We should only choose (3) if we really need to. Often (2) will be sufficient.

# Options if we want a modelling approach

- ② Fit model separately in each country
  - Save model object
  - Send to hub



# Example

- Detailed example on webpage [pclambert.net/XXX](http://pclambert.net/XXX)
- Uses entirely simulated (synthetic) data, so code and data available for people to try for themselves.
- Comparing Country A and Country B.
- Assume I do not have access to data in Country B.

- I have a colleague willing to run code in Country A

## Fit model in Country A

```
// Fit relative survival model
. stpm3 i.sex##@ns(age,df(3)), scale(lncumhazard) df(3) ///
>          bhazard(rate) tvc(i.sex @ns(age,df(3)))
// Save model object
. estimates save countryA.ster
```

- I have a colleague willing to run code in Country A

## Fit model in Country A

```
// Fit relative survival model
. stpm3 i.sex##@ns(age,df(3)), scale(lncumhazard) df(3) ///
>          bhazard(rate) tvc(i.sex @ns(age,df(3)))
// Save model object
. estimates save countryA.ster
```

- My colleague sends this file me in Country B (or elsewhere)

# What's stored in .ster file?

- The ingredients needed to predict survival etc from the model.
  - Names of covariates included in the model
  - Parameter estimates and variances
  - Knot locations for spline functions.
  - Various other details (Number of parameters, sample size, likelihood etc)

# What's stored in .ster file?

- The ingredients needed to predict survival etc from the model.
  - Names of covariates included in the model
  - Parameter estimates and variances
  - Knot locations for spline functions.
  - Various other details (Number of parameters, sample size, likelihood etc)

Crucially it contains no individual level data

# Analyse data in Country B

```
// Load data
. use https://www.pclambert.net/data/CountryB, clear

// Fit relative survival model
. stpm3 i.sex##@ns(age,df(3)), scale(lncumhazard) df(3) ///
>          bhazard(rate) tvc(i.sex @ns(age,df(3)))

// standardized relative survival for Country B
. standsurv RS_B, surv timevar(tt) ci frame(RS, replace)

// now load model object for Country A (BUT NOT DATA)
. estimate use CountryA

// standardized relative survival for Country B
// NOTE: standardized to age/sex distribution of Country A
. standsurv RS_A, surv frame(RS, merge) ci
```

# Analyse data in Country B

```
// Load data
. use https://www.pclambert.net/data/CountryB, clear

// Fit relative survival model
. stpm3 i.sex##@ns(age,df(3)), scale(lncumhazard) df(3) ///
>          bhazard(rate) tvc(i.sex @ns(age,df(3)))

// standardized relative survival for Country B
. standsurv RS_B, surv timevar(tt) ci frame(RS, replace)

// now load model object for Country A (BUT NOT DATA)
. estimate use CountryA

// standardized relative survival for Country B
// NOTE: standardized to age/sex distribution of Country A
. standsurv RS_A, surv frame(RS, merge) ci
```

# Analyse data in Country B

```
// Load data
. use https://www.pclambert.net/data/CountryB, clear

// Fit relative survival model
. stpm3 i.sex##@ns(age,df(3)), scale(lncumhazard) df(3) ///
>          bhazard(rate) tvc(i.sex @ns(age,df(3)))

// standardized relative survival for Country B
. standsurv RS_B, surv timevar(tt) ci frame(RS, replace)

// now load model object for Country A (BUT NOT DATA)
. estimate use CountryA

// standardized relative survival for Country B
// NOTE: standardized to age/sex distribution of Country A
. standsurv RS_A, surv frame(RS, merge) ci
```



# Analyse data in Country B

```
// Load data
. use https://www.pclambert.net/data/CountryB, clear

// Fit relative survival model
. stpm3 i.sex##@ns(age,df(3)), scale(lncumhazard) df(3) ///
>          bhazard(rate) tvc(i.sex @ns(age,df(3)))

// standardized relative survival for Country B
. standsurv RS_B, surv timevar(tt) ci frame(RS, replace)

// Load model object for Country A (BUT NOT DATA)
. estimate use CountryA

// standardized relative survival for Country B
// NOTE: standardized to age/sex distribution of Country A
. standsurv RS_A, surv frame(RS, merge) ci
```

# Analyse data in Country B

```
// Load data
. use https://www.pclambert.net/data/CountryB, clear

// Fit relative survival model
. stpm3 i.sex##@ns(age,df(3)), scale(lncumhazard) df(3) ///
>          bhazard(rate) tvc(i.sex @ns(age,df(3)))

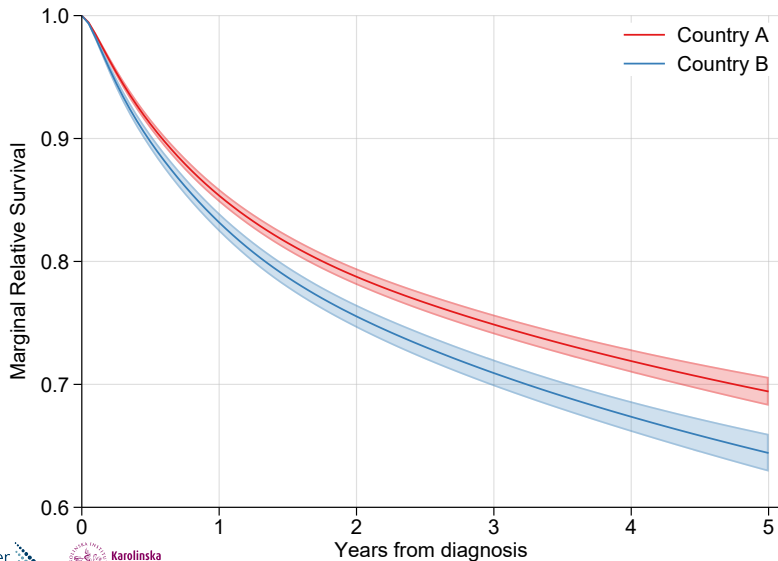
// standardized relative survival for Country B
. standsurv RS_B, surv timevar(tt) ci frame(RS, replace)

// now load model object for Country A (BUT NOT DATA)
. estimate use CountryA

// standardized relative survival for Country B
// NOTE: standardized to age/sex distribution of Country A
. standsurv RS_A, surv frame(RS, merge) ci
```

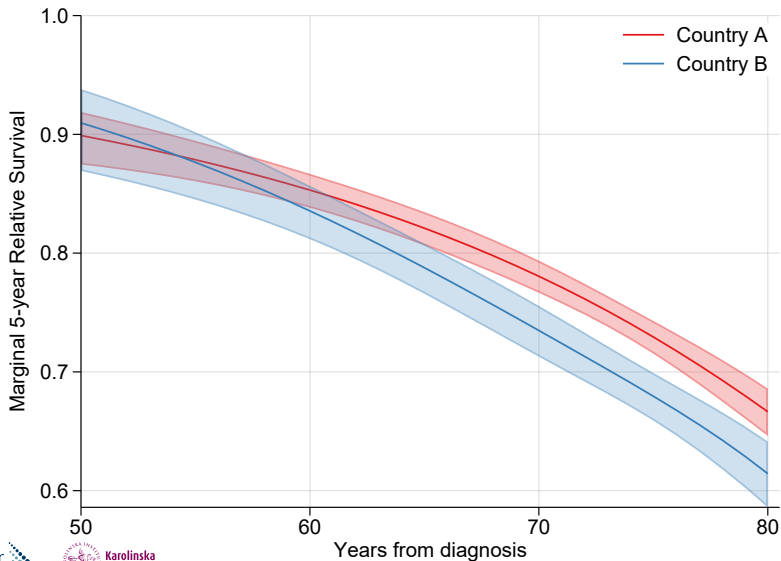
# Results (of simulated data)

Age standardized relative survival



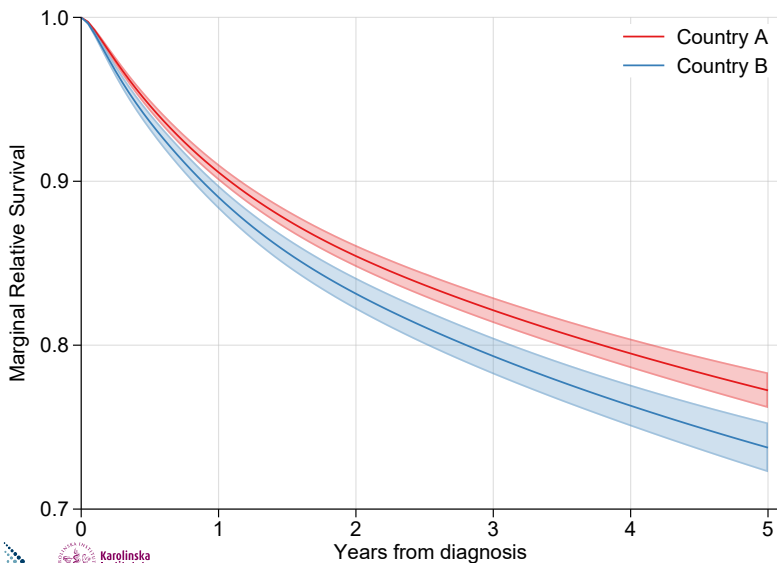
# Results (of simulated data)

5 year relative survival as a function of age

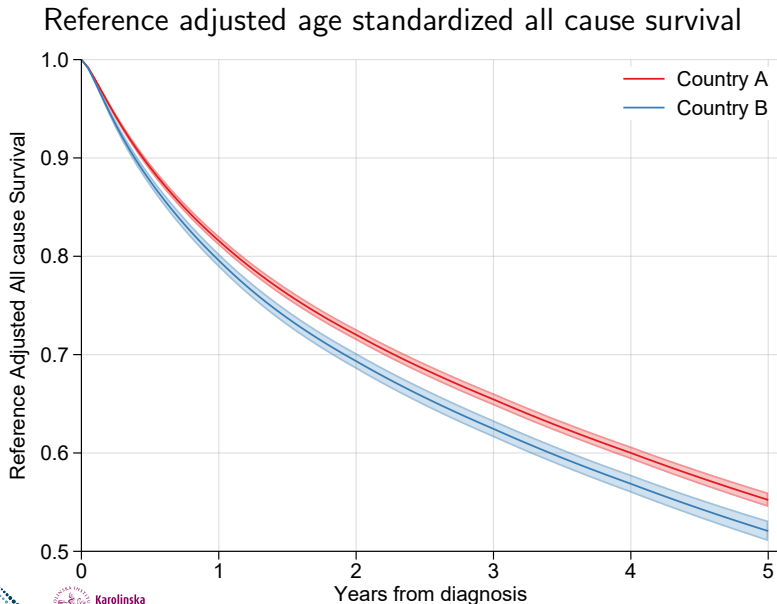


# Results (of simulated data)

Age standardized relative survival (to ICSS age groups)



# Results (of simulated data)



# Standardization

- In my example standardization was to the age/sex distribution of Country B
- Easy to standardize to to external reference, e.g ICSS.
- Also possible to standardize to age/sex distribution of Country A with some summary (aggregated) information.

[See tutorial for examples](#)

# Discussion

- Simple way to fit statistics models, but still obtain useful, and comparable, summaries from those models.
- More flexible than each country producing summaries and just sending those.
- Data quality, inclusion/exclusion criteria, consistency of variable naming/labelling very important.
- Not have the control and ability to fit a combined model, but is an easy practical solution that works.
- Brief overview, lots more details on my webpage.