



Faculty of Health Sciences



# $R^2$ -type Curves for Dynamic Predictions from Joint Longitudinal-Survival Models

Inference & application to prediction of kidney graft failure

Paul Blanche

joint work with

M-C. Fournier & E. Dantan (Nantes, France)

July 2015  
Slide 1/29



## Context & Motivation

- Medical researchers hope to improve patient management using **earlier diagnoses**
- Statisticians can help by fitting **prediction models**
- The making of so-called "**dynamic**" predictions has recently received a lot of attention
- In order to be useful for medical practice, predictions should be "accurate"

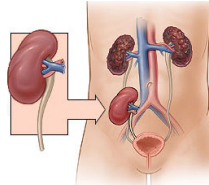
*How should we evaluate dynamic prediction accuracy?*



# Data & Clinical goal

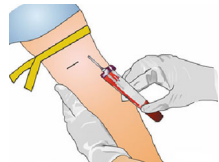
## ► Data:

DIVAT cohort data  
of kidney transplant recipients  
(subsample,  $n = 4,119$ )



## ► Clinical goal:

- Dynamic prediction of risk of kidney graft failure (death or return to dialysis)
- Using repeated measurements of serum creatinine



# DIVAT data sample (n=4,119)

- French cohort
- Adult recipients
- Transplanted after 2000
- Creatinine measured every year

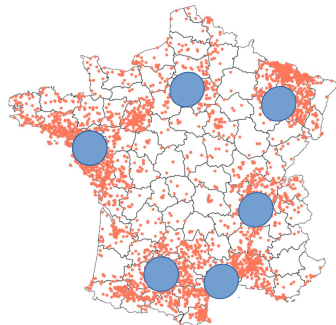


([www.divat.fr](http://www.divat.fr))



# DIVAT data sample (n=4,119)

- French cohort
- Adult recipients
- Transplanted after 2000
- Creatinine measured every year
- 6 centers



([www.divat.fr](http://www.divat.fr))



# Statistical challenges discussed

*How to **evaluate** and/or compare dynamic predictions?*

► **Using concepts** of:

- Discrimination
- Calibration

► **Accounting** for:

- Dynamic setting
- Censoring issue



# Basic idea & concepts for evaluating predictions

*Basic idea: comparing predictions and observations*



# Basic idea & concepts for evaluating predictions

*Basic idea: comparing predictions and observations (simple!)*





# Basic idea & concepts for evaluating predictions

*Basic idea: comparing predictions and observations (simple!)*

## Concepts:

▶ **Discrimination:**

▶ **Calibration:**



# Basic idea & concepts for evaluating predictions

*Basic idea: comparing predictions and observations (simple!)*

## Concepts:

### ► **Discrimination:**

A model has high discriminative power if the range of predicted risks is wide and subjects with low (high) predicted risk are more (less) likely to experience the event.

### ► **Calibration:**



# Basic idea & concepts for evaluating predictions

*Basic idea: comparing predictions and observations (simple!)*

## Concepts:

### ► **Discrimination:**

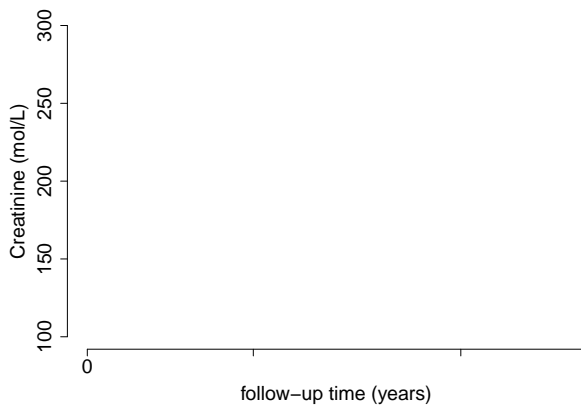
A model has high discriminative power if the range of predicted risks is wide and subjects with low (high) predicted risk are more (less) likely to experience the event.

### ► **Calibration:**

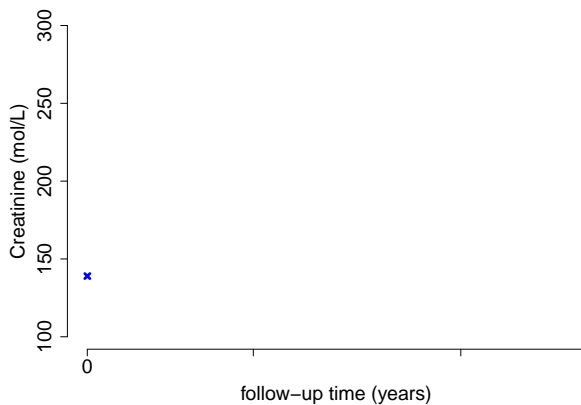
A model is calibrated if we can expect that  $x$  subjects out of 100 experience the event among all subjects that receive a predicted risk of  $x\%$  (“**weak**” definition).



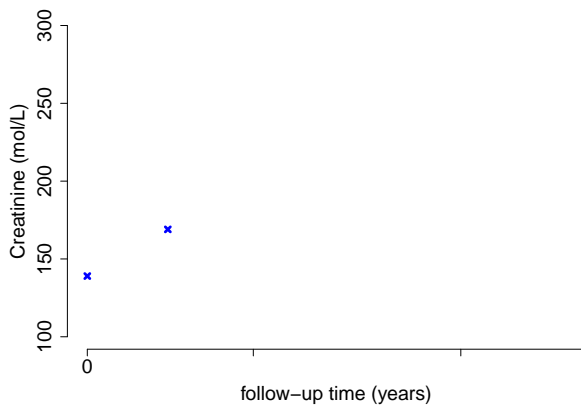
# Dynamic prediction



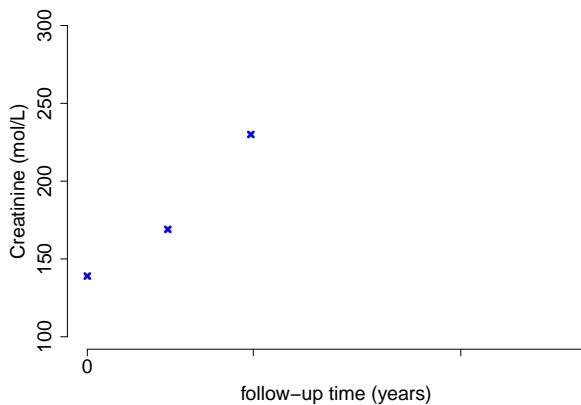
# Dynamic prediction



# Dynamic prediction

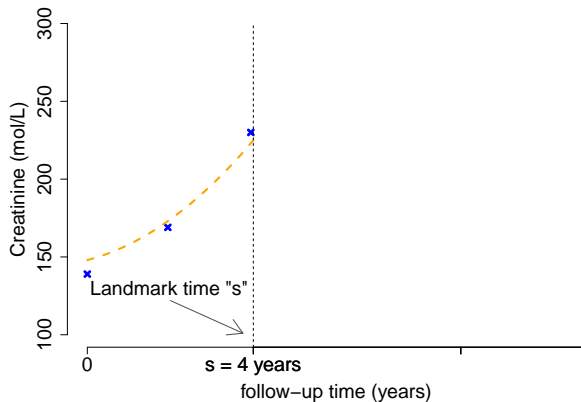


# Dynamic prediction



# Dynamic prediction

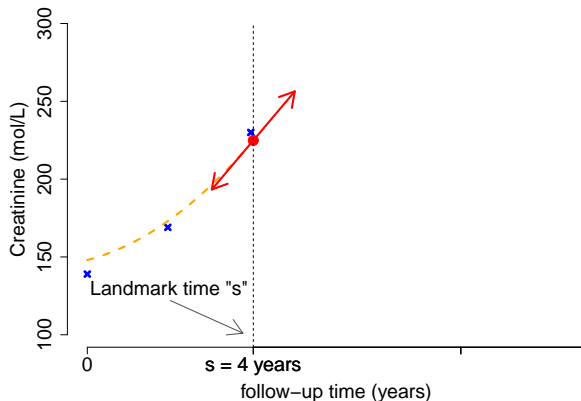
- $s$ : Landmark time at which predictions are made (varies)





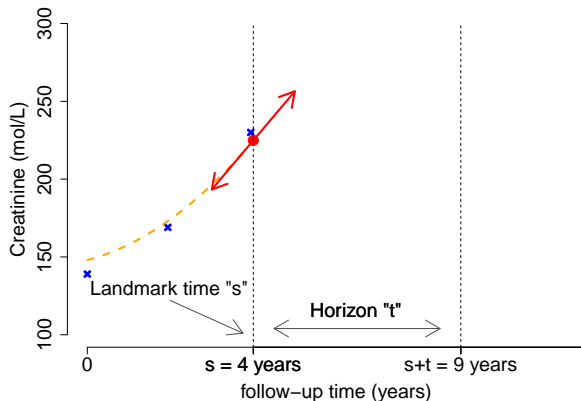
# Dynamic prediction

- $s$ : Landmark time at which predictions are made (varies)



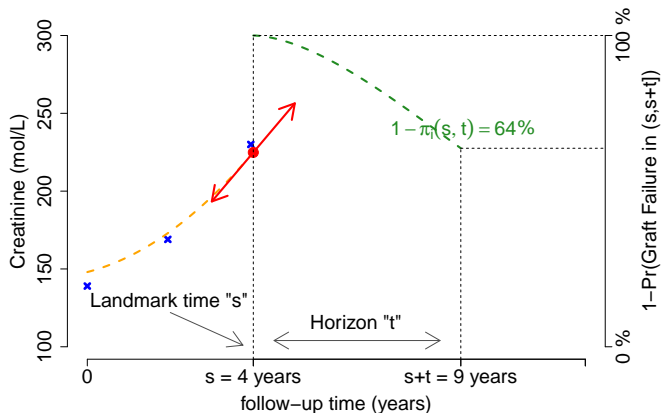
# Dynamic prediction

- $s$ : Landmark time at which predictions are made (**varies**)
- $t$ : prediction horizon (**fixed**)

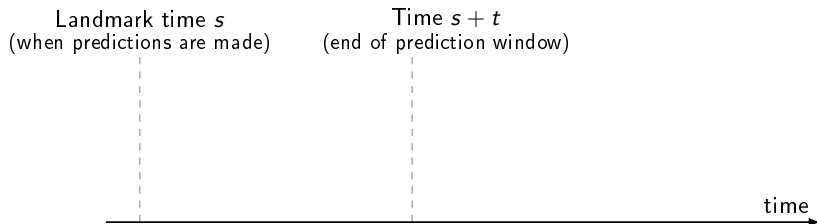


# Dynamic prediction

- $s$ : Landmark time at which predictions are made (varies)
- $t$ : prediction horizon (fixed)



# Right censoring issue

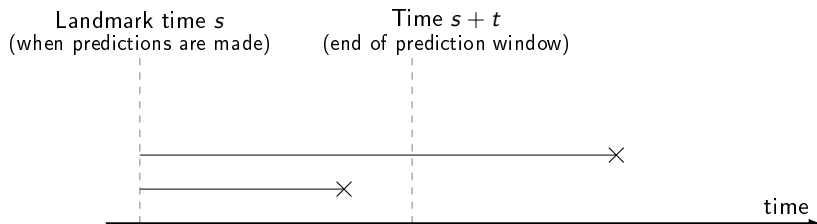


$$D_i(s, t) = \mathbb{1}\{\text{event occurs in } (s, s + t]\}$$



# Right censoring issue

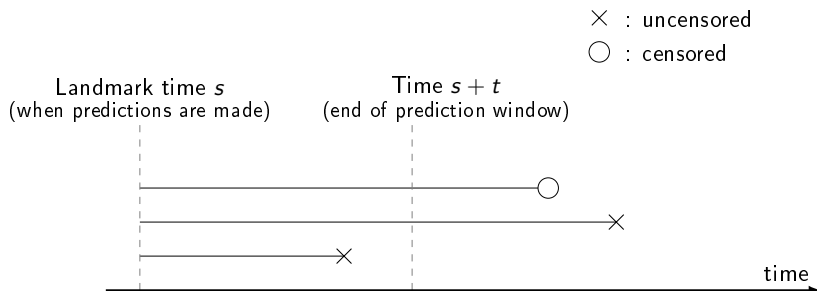
× : uncensored



$$D_i(s, t) = \mathbb{1}\{\text{event occurs in } (s, s + t]\}$$



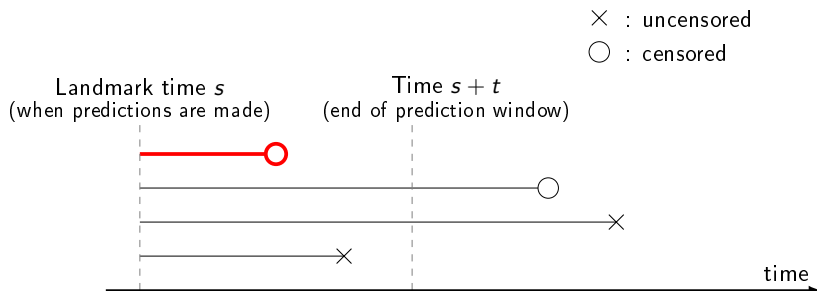
# Right censoring issue



$$D_i(s, t) = \mathbb{1}\{\text{event occurs in } (s, s + t]\}$$



## Right censoring issue



For subject  $i$  censored within  $(s, s + t]$  the **status**

$$D_i(s, t) = \mathbb{1}\{\text{event occurs in } (s, s + t]\}$$

is unknown.



# Notations for population parameters

► **Indicator of event** in  $(s, s + t]$ :

$$D_i(s, t) = \mathbb{1}\{s < T_i \leq s + t\}$$

where  $T_i$  is the time-to-event.





# Notations for population parameters

- ▶ **Indicator of event** in  $(s, s + t]$ :

$$D_i(s, t) = \mathbb{1}\{s < T_i \leq s + t\}$$

where  $T_i$  is the time-to-event.

- ▶ **Dynamic predictions:**

$$\pi_i(s, t)$$



# Notations for population parameters

- **Indicator of event** in  $(s, s + t]$ :

$$D_i(s, t) = \mathbb{1}\{s < T_i \leq s + t\}$$

where  $T_i$  is the time-to-event.

- **Dynamic predictions:**

$$\pi_i(s, t) = \widehat{\mathbb{P}}_{\widehat{\xi}} \left( \quad \right)$$

- $\widehat{\xi}$ : previously estimated parameters (from independent training data)



# Notations for population parameters

- **Indicator of event** in  $(s, s + t]$ :

$$D_i(s, t) = \mathbb{1}\{s < T_i \leq s + t\}$$

where  $T_i$  is the time-to-event.

- **Dynamic predictions:**

$$\pi_i(s, t) = \widehat{\mathbb{P}}_{\widehat{\xi}} \left( D_i(s, t) = 1 \mid \right)$$

- $\widehat{\xi}$ : previously estimated parameters (from independent training data)



# Notations for population parameters

- **Indicator of event** in  $(s, s + t]$ :

$$D_i(s, t) = \mathbb{1}\{s < T_i \leq s + t\}$$

where  $T_i$  is the time-to-event.

- **Dynamic predictions:**

$$\pi_i(s, t) = \widehat{\mathbb{P}}_{\widehat{\xi}} \left( D_i(s, t) = 1 \mid T_i > s, \mathcal{Y}_i(s), \right)$$

- $\widehat{\xi}$ : previously estimated parameters (from independent training data)
- $\mathcal{Y}_i(s)$ : marker measurements observed before time  $s$



# Notations for population parameters

- **Indicator of event** in  $(s, s + t]$ :

$$D_i(s, t) = \mathbb{1}\{s < T_i \leq s + t\}$$

where  $T_i$  is the time-to-event.

- **Dynamic predictions:**

$$\pi_i(s, t) = \widehat{\mathbb{P}}_{\widehat{\xi}} \left( D_i(s, t) = 1 \mid T_i > s, \mathcal{Y}_i(s), \mathbf{X}_i \right)$$

- $\widehat{\xi}$ : previously estimated parameters (from independent training data)
- $\mathcal{Y}_i(s)$ : marker measurements observed before time  $s$
- $\mathbf{X}_i$ : baseline covariates



## Predictive accuracy

*How close are the predicted risks  $\pi_i(s, t)$  to the “true underlying” risk  $\mathbb{P}(\text{event occurs in } (s, s+t] | \text{information at } s)$ ?*



## Predictive accuracy

*How close are the predicted risks  $\pi_i(s, t)$  to the “true underlying” risk  $\mathbb{P}(\text{event occurs in } (s, s+t] | \text{information at } s)$ ?*

► **Prediction Error:**

$$\text{PE}_\pi(s, t) = \mathbb{E} \left[ \left\{ D(s, t) - \pi(s, t) \right\}^2 \mid T > s \right]$$



## Predictive accuracy

*How close are the predicted risks  $\pi_i(s, t)$  to the “true underlying” risk  $\mathbb{P}(\text{event occurs in } (s, s+t] | \text{information at } s)$ ?*

► **Prediction Error:**

$$\text{PE}_\pi(s, t) = \mathbb{E} \left[ \left\{ D(s, t) - \pi(s, t) \right\}^2 \mid T > s \right]$$

- the lower the better
- $\text{PE} \approx \text{Bias}^2 + \text{Variance}$
- evaluates both **Calibration** and **Discrimination**
- depends on  $\mathbb{P}(\text{event occurs in } (s, s+t] | \text{at risk at } s)$
- often called "Expected **Brier Score**"





How does the PE relate to calibration and discrimination?

$$PE_{\pi}(s, t) = \mathbb{E} \left[ \left\{ D(s, t) - \pi(s, t) \right\}^2 \mid T > s \right]$$



How does the PE relate to calibration and discrimination?

$$\text{PE}_\pi(s, t) = \mathbb{E} \left[ \left\{ D(s, t) - \mathbb{E}[D(s, t) | \mathcal{H}^\pi(s)] \right. \right. \\ \left. \left. + \underbrace{\mathbb{E}[D(s, t) | \mathcal{H}^\pi(s)]}_{\text{"true underlying" risk}} - \pi(s, t) \right\}^2 \middle| T > s \right]$$



## How does the PE relate to calibration and discrimination?

$$\begin{aligned} \text{PE}_\pi(s, t) = & \mathbb{E} \left[ \left\{ D(s, t) - \mathbb{E}[D(s, t) | \mathcal{H}^\pi(s)] \right\}^2 \middle| T > s \right] \\ & + \mathbb{E} \left[ \left\{ \underbrace{\mathbb{E}[D(s, t) | \mathcal{H}^\pi(s)]}_{\text{"true underlying" risk}} - \pi(s, t) \right\}^2 \middle| T > s \right] \end{aligned}$$

$\mathcal{H}^\pi(s) = \{\mathcal{X}^\pi(s), T > s\}$  denotes the subject-specific history at time  $s$ .



# How does the PE relate to calibration and discrimination?

$$\begin{aligned}
 \text{PE}_\pi(s, t) = & \underbrace{\mathbb{E} \left[ \left\{ D(s, t) - \mathbb{E}[D(s, t) | \mathcal{H}^\pi(s)] \right\}^2 \middle| T > s \right]}_{\text{Inseparability}} \\
 & + \underbrace{\mathbb{E} \left[ \left\{ \mathbb{E}[D(s, t) | \mathcal{H}^\pi(s)] - \pi(s, t) \right\}^2 \middle| T > s \right]}_{\text{Bias/Calibration}}
 \end{aligned}$$

$\mathcal{H}^\pi(s) = \{\mathcal{X}^\pi(s), T > s\}$  denotes the subject-specific history at time  $s$ .



# How does the PE relate to calibration and discrimination?

$$\begin{aligned}
 \text{PE}_\pi(s, t) = & \mathbb{E} \left[ \underbrace{\text{Var}\{D(s, t) | \mathcal{H}^\pi(s)\}}_{\text{Discrimination}} \middle| T > s \right] \\
 & + \mathbb{E} \left[ \underbrace{\left\{ \mathbb{E}[D(s, t) | \mathcal{H}^\pi(s)] - \pi(s, t) \right\}^2}_{\text{Calibration}} \middle| T > s \right]
 \end{aligned}$$

$\mathcal{H}^\pi(s) = \{\mathcal{X}^\pi(s), T > s\}$  denotes the subject-specific history at time  $s$ .



# How does the PE relate to calibration and discrimination?

$$\begin{aligned}
 \text{PE}_\pi(s, t) = & \underbrace{\mathbb{E} \left[ \text{Var} \{ D(s, t) | \mathcal{H}^\pi(s) \} \mid T > s \right]}_{\text{Discrimination}} \\
 & + \underbrace{\mathbb{E} \left[ \left\{ \mathbb{E} [ D(s, t) | \mathcal{H}^\pi(s) ] - \pi(s, t) \right\}^2 \mid T > s \right]}_{\text{Calibration}}
 \end{aligned}$$

$\mathcal{H}^\pi(s) = \{ \mathcal{X}^\pi(s), T > s \}$  denotes the subject-specific history at time  $s$ .

- the more discriminating  $\mathcal{H}^\pi(s)$  the smaller  $\text{Var} \{ D(s, t) | \mathcal{H}^\pi(s) \}$



# How does the PE relate to calibration and discrimination?

$$\begin{aligned}
 \text{PE}_\pi(s, t) = & \underbrace{\mathbb{E} \left[ \text{Var} \{ D(s, t) | \mathcal{H}^\pi(s) \} \mid T > s \right]}_{\text{Discrimination}} \\
 & + \underbrace{\mathbb{E} \left[ \left\{ \mathbb{E} [ D(s, t) | \mathcal{H}^\pi(s) ] - \pi(s, t) \right\}^2 \mid T > s \right]}_{\text{Calibration}}
 \end{aligned}$$

$\mathcal{H}^\pi(s) = \{ \mathcal{X}^\pi(s), T > s \}$  denotes the subject-specific history at time  $s$ .

- ▶ the more discriminating  $\mathcal{H}^\pi(s)$  the smaller  $\text{Var} \{ D(s, t) | \mathcal{H}^\pi(s) \}$
- ▶  $\mathbb{E} [ D(s, t) | \mathcal{H}^\pi(s) ] - \pi(s, t) \equiv 0$  defines "strong" calibration.



## How does the PE relate to calibration and discrimination?

$$\begin{aligned}
 \text{PE}_\pi(s, t) = & \underbrace{\mathbb{E} \left[ \mathbf{Var} \{ D(s, t) | \mathcal{H}^\pi(s) \} \mid T > s \right]}_{\text{Does NOT depend on } \pi(s, t)} \\
 & + \underbrace{\mathbb{E} \left[ \left\{ \mathbb{E} [ D(s, t) | \mathcal{H}^\pi(s) ] - \pi(s, t) \right\}^2 \mid T > s \right]}_{\text{Depends on } \pi(s, t)}
 \end{aligned}$$

$\mathcal{H}^\pi(s) = \{ \mathcal{X}^\pi(s), T > s \}$  denotes the subject-specific history at time  $s$ .

- ▶ the more discriminating  $\mathcal{H}^\pi(s)$  the smaller  $\mathbf{Var} \{ D(s, t) | \mathcal{H}^\pi(s) \}$
- ▶  $\mathbb{E} [ D(s, t) | \mathcal{H}^\pi(s) ] - \pi(s, t) \equiv 0$  defines "strong" calibration.





## $R_{\pi}^2$ -type criterion

### ► Benchmark $PE_0$

The **best "null" prediction tool**, which gives the same (marginal) predicted risk

$$S(s + t|s) = \mathbb{E}[D(s, t)|\mathcal{H}^0(s)], \quad \mathcal{H}^0(s) = \{T > s\}$$

to all subjects leads to

$$PE_0(s, t) = \text{Var}\{D(s, t)|T > s\} = S(s + t|s)\{1 - S(s + t|s)\}.$$



## $R_\pi^2$ -type criterion

### ► Benchmark $PE_0$

The **best "null" prediction tool**, which gives the same (marginal) predicted risk

$$S(s+t|s) = \mathbb{E}[D(s,t)|\mathcal{H}^0(s)], \quad \mathcal{H}^0(s) = \{T > s\}$$

to all subjects leads to

$$PE_0(s,t) = \text{Var}\{D(s,t)|T > s\} = S(s+t|s)\{1 - S(s+t|s)\}.$$

### ► Simple idea

$$R_\pi^2(s,t) = 1 - \frac{PE_\pi(s,t)}{PE_0(s,t)}$$



## Why bother?

- ▶  $R_{\pi}^2(s, t)$  aims to circumvent the difficult interpretation of:
  - the scale on which  $PE(s, t)$  is measured
  - interpretation for trend of  $PE(s, t)$  vs  $s$



## Why bother?

- ▶  $R_{\pi}^2(s, t)$  aims to circumvent the difficult interpretation of:
  - the scale on which  $PE(s, t)$  is measured
  - interpretation for trend of  $PE(s, t)$  vs  $s$
  
- ▶ Because the meaning of **the scale on which  $PE(s, t)$  is measured changes with  $s$** , an increasing/decreasing trend can be due to changes in:
  - the quality of the predictions**and/or**
  - the at risk population



# Changes in the quality of the predictions

*"Essentially, all models are wrong, but some are useful."*, G. Box



# Changes in the quality of the predictions

*"Essentially, all models are wrong, but some are useful."*, G. Box



► The prediction model from which we have obtained the predictions can be “more wrong” for some  $s$  than for some others.

- Calibration term of  $PE(s, t)$  changes with  $s$
- We can work on it!



# Changes in the at risk population

## An example:

- *Patients with cardiovascular history (CV) all die early.*
- *Only those without CV remain at risk for late  $s$ .*
- *Then:*
  - *the earlier  $s$  the more homogeneous the at risk population*
  - *CV is useful for prediction for early  $s$  but useless for late  $s$ .*



# Changes in the at risk population

## An example:

- *Patients with cardiovascular history (CV) all die early.*
- *Only those without CV remain at risk for late  $s$ .*
- *Then:*
  - *the earlier  $s$  the more homogeneous the at risk population*
  - *CV is useful for prediction for early  $s$  but useless for late  $s$ .*

► The available information can be more informative for some  $s$  than for some others.

- Discriminating term of  $PE(s, t)$  changes with  $s$
- This is just how it is, there is nothing we can do!

(we can only work with the data we have)





## $R_{\pi}^2(s, t)$ interpretation

► **Always true:**

Measure of how the prediction tool  $\pi(s, t)$  performs compared to the benchmark null prediction tool, which gives the same predicted risk to all subjects (marginal risk).



## $R_{\pi}^2(s, t)$ interpretation

► **Always true:**

Measure of how the prediction tool  $\pi(s, t)$  performs compared to the benchmark null prediction tool, which gives the same predicted risk to all subjects (marginal risk).

► **When** predictions are **calibrated** (strongly):

$$R_{\pi}^2(s, t) = \frac{\text{Var}\{\pi(s, t) | T > s\}}{\text{Var}\{D(s, t) | T > s\}}$$

**explained variation**

(after little algebra)



## $R_{\pi}^2(s, t)$ interpretation

► **Always true:**

Measure of how the prediction tool  $\pi(s, t)$  performs compared to the benchmark null prediction tool, which gives the same predicted risk to all subjects (marginal risk).

► **When** predictions are **calibrated** (strongly):

$$R_{\pi}^2(s, t) = \frac{\text{Var}\{\pi(s, t) | T > s\}}{\text{Var}\{D(s, t) | T > s\}} \quad \text{explained variation}$$

$$= \text{Corr}^2\{D(s, t), \pi(s, t) | T > s\} \quad \text{correlation}$$

(after little algebra)



## $R_{\pi}^2(s, t)$ interpretation

► **Always true:**

Measure of how the prediction tool  $\pi(s, t)$  performs compared to the benchmark null prediction tool, which gives the same predicted risk to all subjects (marginal risk).

► **When** predictions are **calibrated** (strongly):

$$R_{\pi}^2(s, t) = \frac{\text{Var}\{\pi(s, t) | T > s\}}{\text{Var}\{D(s, t) | T > s\}} \quad \text{explained variation}$$

$$= \text{Corr}^2\{D(s, t), \pi(s, t) | T > s\} \quad \text{correlation}$$

$$= \mathbb{E}\left\{\pi(s, t) \mid D(s, t) = 1, T > s\right\} - \mathbb{E}\left\{\pi(s, t) \mid D(s, t) = 0, T > s\right\} \quad \text{mean risk difference}$$

(after little algebra)



# Observations & IPCW PE estimator

## ► **Observations** (i.i.d.)

$$\left\{ (\tilde{T}_i, \Delta_i, \pi_i(\cdot, \cdot)), i = 1, \dots, n \right\} \quad \text{where } \tilde{T}_i = T_i \wedge C_i, \Delta_i = \mathbb{1}\{T_i \leq C_i\}$$



## Observations & IPCW PE estimator

► **Observations** (i.i.d.)

$$\left\{ (\tilde{T}_i, \Delta_i, \pi_i(\cdot, \cdot)), i = 1, \dots, n \right\} \quad \text{where } \tilde{T}_i = T_i \wedge C_i, \Delta_i = \mathbb{1}\{T_i \leq C_i\}$$

► Indicator of “**observed event occurrence**” in  $(s, s + t]$ :

$$\tilde{D}_i(s, t) = \mathbb{1}\{s < \tilde{T}_i \leq s + t, \Delta_i = 1\} = \begin{cases} 1 & : \text{event occurred} \\ 0 & : \text{event did not occur} \\ & \text{or censored obs.} \end{cases}$$



## Observations & IPCW PE estimator

► **Observations** (i.i.d.)

$$\left\{ (\tilde{T}_i, \Delta_i, \pi_i(\cdot, \cdot)), i = 1, \dots, n \right\} \quad \text{where } \tilde{T}_i = T_i \wedge C_i, \Delta_i = \mathbb{1}\{T_i \leq C_i\}$$

► Indicator of “**observed event occurrence**” in  $(s, s + t]$ :

$$\tilde{D}_i(s, t) = \mathbb{1}\{s < \tilde{T}_i \leq s + t, \Delta_i = 1\} = \begin{cases} 1 & \text{: event occurred} \\ 0 & \text{: event did not occur} \\ & \text{or censored obs.} \end{cases}$$

► Inverse Probability of Censoring Weighting (IPCW) estimator:

$$\widehat{PE}_\pi(s, t) = \frac{1}{n} \sum_{i=1}^n \left\{ \tilde{D}_i(s, t) - \pi_i(s, t) \right\}^2$$



## Observations & IPCW PE estimator

► **Observations** (i.i.d.)

$$\left\{ (\tilde{T}_i, \Delta_i, \pi_i(\cdot, \cdot)), i = 1, \dots, n \right\} \quad \text{where } \tilde{T}_i = T_i \wedge C_i, \Delta_i = \mathbb{1}\{T_i \leq C_i\}$$

► Indicator of “**observed event occurrence**” in  $(s, s + t]$ :

$$\tilde{D}_i(s, t) = \mathbb{1}\{s < \tilde{T}_i \leq s + t, \Delta_i = 1\} = \begin{cases} 1 & : \text{event occurred} \\ 0 & : \text{event did not occur} \\ & \text{or censored obs.} \end{cases}$$

► Inverse Probability of Censoring Weighting (IPCW) estimator:

$$\widehat{PE}_\pi(s, t) = \frac{1}{n} \sum_{i=1}^n \widehat{W}_i(s, t) \left\{ \tilde{D}_i(s, t) - \pi_i(s, t) \right\}^2$$





## Observations & IPCW PE estimator

### ► Observations (i.i.d.)

$$\left\{ (\tilde{T}_i, \Delta_i, \pi_i(\cdot, \cdot)), i = 1, \dots, n \right\} \quad \text{where } \tilde{T}_i = T_i \wedge C_i, \Delta_i = \mathbb{1}\{T_i \leq C_i\}$$

### ► Indicator of “observed event occurrence” in $(s, s + t]$ :

$$\tilde{D}_i(s, t) = \mathbb{1}\{s < \tilde{T}_i \leq s + t, \Delta_i = 1\} = \begin{cases} 1 & \text{: event occurred} \\ 0 & \text{: event did not occur} \\ & \text{or censored obs.} \end{cases}$$

### ► Inverse Probability of Censoring Weighting (IPCW) estimator:

$$\widehat{PE}_\pi(s, t) = \frac{1}{n} \sum_{i=1}^n \widehat{W}_i(s, t) \left\{ \tilde{D}_i(s, t) - \pi_i(s, t) \right\}^2$$

and

$$\widehat{R}_\pi^2(s, t) = 1 - \frac{\widehat{PE}_\pi(s, t)}{\widehat{PE}_0(s, t)}$$



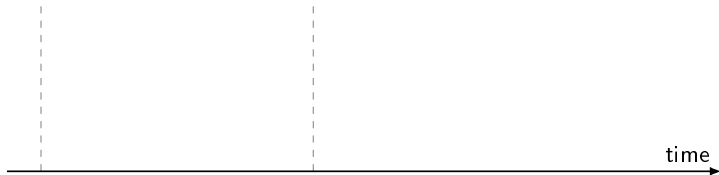
# Inverse Probability of Censoring Weights

$$\widehat{W}_i(s, t) = \quad + \quad +$$

with  $\widehat{G}(u|s)$  the Kaplan-Meier estimator of  $\mathbb{P}(C > u | C > s)$ .

Landmark time  $s$

Time  $s + t$



# Inverse Probability of Censoring Weights

$$\widehat{W}_i(s, t) = \frac{\mathbb{1}\{s < \widetilde{T}_i \leq s + t\} \Delta_i}{\widehat{G}(\widetilde{T}_i|s)} + \quad +$$

with  $\widehat{G}(u|s)$  the Kaplan-Meier estimator of  $\mathbb{P}(C > u|C > s)$ .

Landmark time  $s$

Time  $s + t$

× : uncensored

○ : censored



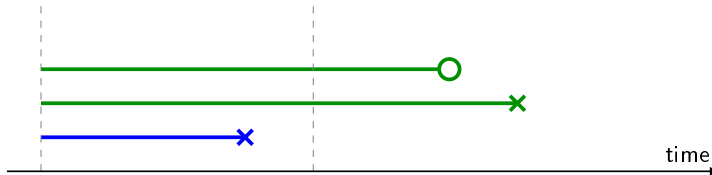
# Inverse Probability of Censoring Weights

$$\widehat{W}_i(s, t) = \frac{\mathbb{1}\{s < \widetilde{T}_i \leq s + t\} \Delta_i}{\widehat{G}(\widetilde{T}_i|s)} + \frac{\mathbb{1}\{\widetilde{T}_i > s + t\}}{\widehat{G}(s + t|s)} +$$

with  $\widehat{G}(u|s)$  the Kaplan-Meier estimator of  $\mathbb{P}(C > u|C > s)$ .

Landmark time  $s$

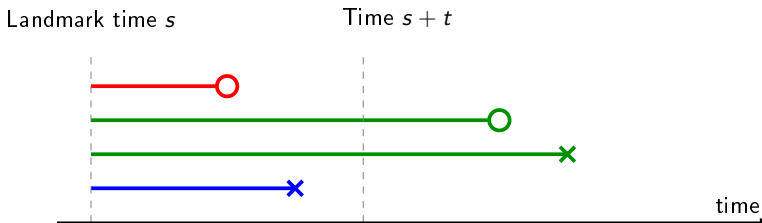
Time  $s + t$



# Inverse Probability of Censoring Weights

$$\widehat{W}_i(s, t) = \frac{\mathbb{1}\{s < \widetilde{T}_i \leq s + t\} \Delta_i}{\widehat{G}(\widetilde{T}_i|s)} + \frac{\mathbb{1}\{\widetilde{T}_i > s + t\}}{\widehat{G}(s + t|s)} + 0$$

with  $\widehat{G}(u|s)$  the Kaplan-Meier estimator of  $\mathbb{P}(C > u|C > s)$ .



## Asymptotic i.i.d. representation

**Lemma:** Assume that the censoring time  $C$  is independent of  $(T, \eta, \pi(\cdot, \cdot))$  and let  $\theta$  denote either  $PE_\pi$ ,  $R_\pi^2$  or a difference in PE or  $R_\pi^2$ , then

$$\sqrt{n} \left( \widehat{\theta}(s, t) - \theta(s, t) \right) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \text{IF}_\theta(\widetilde{T}_i, \Delta_i, \pi_i(s, t), s, t) + o_p(1)$$

where  $\text{IF}_\theta(\widetilde{T}_i, \Delta_i, \pi_i(s, t), s, t)$  being :

- ▶ zero-mean i.i.d. terms
- ▶ easy to estimate (using Nelson-Aalen & Kaplan-Meier)



## Pointwise confidence interval (**fixed s**)

- Asymptotic normality:

$$\sqrt{n}(\widehat{\theta}(s, t) - \theta(s, t)) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \sigma_{s,t}^2)$$

- 95% confidence interval:

$$\left\{ \widehat{\theta}(s, t) \pm z_{1-\alpha/2} \frac{\widehat{\sigma}_{s,t}}{\sqrt{n}} \right\}$$

where  $z_{1-\alpha/2}$  is the  $1 - \alpha/2$  quantile of  $\mathcal{N}(0, 1)$ .

- Variance estimator:

$$\widehat{\sigma}_{s,t}^2 = \frac{1}{n} \sum_{i=1}^n \left\{ \widehat{\text{IF}}_{\theta}(\widetilde{T}_i, \Delta_i, \pi_i(s, t), s, t) \right\}^2$$



# Simultaneous confidence band over $s \in \mathcal{S}$

$$\left\{ \hat{\theta}(s, t) \pm \hat{q}_{1-\alpha}^{(S, t)} \frac{\hat{\sigma}_{s, t}}{\sqrt{n}} \right\}, \quad s \in \mathcal{S}$$





# Simultaneous confidence band over $s \in \mathcal{S}$

$$\left\{ \hat{\theta}(s, t) \pm \hat{q}_{1-\alpha}^{(\mathcal{S}, t)} \frac{\hat{\sigma}_{s, t}}{\sqrt{n}} \right\}, \quad s \in \mathcal{S}$$

Computation of  $\hat{q}_{1-\alpha}^{(\mathcal{S}, t)}$  by the **simulation algorithm**

- 1 For  $b = 1, \dots, B$ , say  $B = 4000$ , do:



# Simultaneous confidence band over $s \in \mathcal{S}$

$$\left\{ \hat{\theta}(s, t) \pm \hat{q}_{1-\alpha}^{(\mathcal{S}, t)} \frac{\hat{\sigma}_{s, t}}{\sqrt{n}} \right\}, \quad s \in \mathcal{S}$$

Computation of  $\hat{q}_{1-\alpha}^{(\mathcal{S}, t)}$  by the **simulation algorithm** ( $\approx$  Wild Bootstrap):

- 1 For  $b = 1, \dots, B$ , say  $B = 4000$ , do:

*(Conditional multiplier central limit theorem)*



# Simultaneous confidence band over $s \in \mathcal{S}$

$$\left\{ \hat{\theta}(s, t) \pm \hat{q}_{1-\alpha}^{(\mathcal{S}, t)} \frac{\hat{\sigma}_{s, t}}{\sqrt{n}} \right\}, \quad s \in \mathcal{S}$$

Computation of  $\hat{q}_{1-\alpha}^{(\mathcal{S}, t)}$  by the **simulation algorithm** ( $\approx$  Wild Bootstrap):

- 1 For  $b = 1, \dots, B$ , say  $B = 4000$ , do:
  - 1 Generate  $\{\omega_1^b, \dots, \omega_n^b\}$  from  $n$  i.i.d.  $\mathcal{N}(0, 1)$ .

(Conditional multiplier central limit theorem)



# Simultaneous confidence band over $s \in \mathcal{S}$

$$\left\{ \hat{\theta}(s, t) \pm \hat{q}_{1-\alpha}^{(S,t)} \frac{\hat{\sigma}_{s,t}}{\sqrt{n}} \right\}, \quad s \in \mathcal{S}$$

Computation of  $\hat{q}_{1-\alpha}^{(S,t)}$  by the **simulation algorithm** ( $\approx$  Wild Bootstrap):

- 1 For  $b = 1, \dots, B$ , say  $B = 4000$ , do:
  - 1 Generate  $\{\omega_1^b, \dots, \omega_n^b\}$  from  $n$  i.i.d.  $\mathcal{N}(0, 1)$ .
  - 2 Using the plug-in estimator  $\hat{\text{IF}}_{\theta}(\cdot)$ , compute:

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\hat{\text{IF}}_{\theta}(\tilde{T}_i, \Delta_i, \pi_i(s, t), s, t)}{\hat{\sigma}_{s,t}}$$

(Conditional multiplier central limit theorem)



# Simultaneous confidence band over $s \in \mathcal{S}$

$$\left\{ \hat{\theta}(s, t) \pm \hat{q}_{1-\alpha}^{(S,t)} \frac{\hat{\sigma}_{s,t}}{\sqrt{n}} \right\}, \quad s \in \mathcal{S}$$

Computation of  $\hat{q}_{1-\alpha}^{(S,t)}$  by the **simulation algorithm** ( $\approx$  Wild Bootstrap):

- 1 For  $b = 1, \dots, B$ , say  $B = 4000$ , do:
  - 1 Generate  $\{\omega_1^b, \dots, \omega_n^b\}$  from  $n$  i.i.d.  $\mathcal{N}(0, 1)$ .
  - 2 Using the plug-in estimator  $\hat{\mathbb{F}}_{\theta}(\cdot)$ , compute:

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \omega_i^b \frac{\hat{\mathbb{F}}_{\theta}(\tilde{T}_i, \Delta_i, \pi_i(s, t), s, t)}{\hat{\sigma}_{s,t}}$$

(Conditional multiplier central limit theorem)



# Simultaneous confidence band over $s \in \mathcal{S}$

$$\left\{ \widehat{\theta}(s, t) \pm \widehat{q}_{1-\alpha}^{(\mathcal{S}, t)} \frac{\widehat{\sigma}_{s, t}}{\sqrt{n}} \right\}, \quad s \in \mathcal{S}$$

Computation of  $\widehat{q}_{1-\alpha}^{(\mathcal{S}, t)}$  by the **simulation algorithm** ( $\approx$  Wild Bootstrap):

- 1 For  $b = 1, \dots, B$ , say  $B = 4000$ , do:
  - 1 Generate  $\{\omega_1^b, \dots, \omega_n^b\}$  from  $n$  i.i.d.  $\mathcal{N}(0, 1)$ .
  - 2 Using the plug-in estimator  $\widehat{\text{IF}}_{\theta}(\cdot)$ , compute:

$$\gamma^b = \sup_{s \in \mathcal{S}} \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \omega_i^b \frac{\widehat{\text{IF}}_{\theta}(\widetilde{T}_i, \Delta_i, \pi_i(s, t), s, t)}{\widehat{\sigma}_{s, t}} \right|$$

(Conditional multiplier central limit theorem)



# Simultaneous confidence band over $s \in \mathcal{S}$

$$\left\{ \widehat{\theta}(s, t) \pm \widehat{q}_{1-\alpha}^{(S,t)} \frac{\widehat{\sigma}_{s,t}}{\sqrt{n}} \right\}, \quad s \in \mathcal{S}$$

Computation of  $\widehat{q}_{1-\alpha}^{(S,t)}$  by the **simulation algorithm** ( $\approx$  Wild Bootstrap):

- 1 For  $b = 1, \dots, B$ , say  $B = 4000$ , do:
  - 1 Generate  $\{\omega_1^b, \dots, \omega_n^b\}$  from  $n$  i.i.d.  $\mathcal{N}(0, 1)$ .
  - 2 Using the plug-in estimator  $\widehat{\text{IF}}_{\theta}(\cdot)$ , compute:

$$\Upsilon^b = \sup_{s \in \mathcal{S}} \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \omega_i^b \frac{\widehat{\text{IF}}_{\theta}(\widetilde{T}_i, \Delta_i, \pi_i(s, t), s, t)}{\widehat{\sigma}_{s,t}} \right|$$

- 2 Compute  $\widehat{q}_{1-\alpha}^{(S,t)}$  as the  $100(1 - \alpha)$ th percentile of  $\{\Upsilon^1, \dots, \Upsilon^B\}$

(Conditional multiplier central limit theorem)



## DIVAT sample

- Population based study of kidney recipients (n=4,119)
- Split the data into training (2/3) and validation (1/3) samples





# DIVAT sample

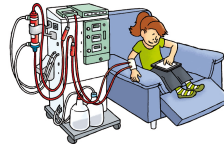
- Population based study of kidney recipients ( $n=4,119$ )
- Split the data into training (2/3) and validation (1/3) samples
- $T$ : time from 1-year after transplantation to graft failure which is:

Death



OR

Return to dialysis

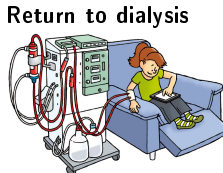


# DIVAT sample

- Population based study of kidney recipients ( $n=4,119$ )
- Split the data into training (2/3) and validation (1/3) samples
- $T$ : time from 1-year after transplantation to graft failure which is:



OR

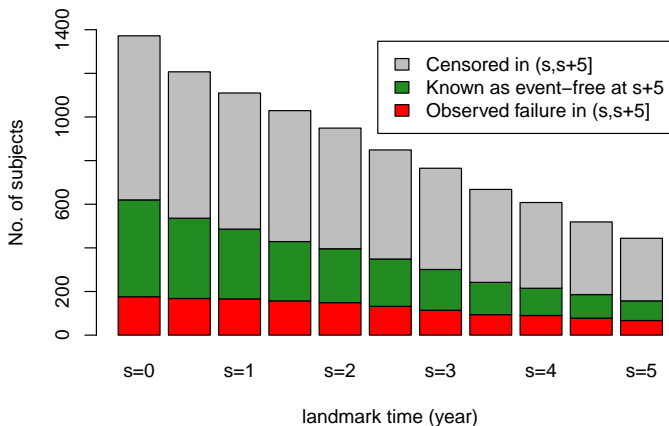


- Censoring due to:
  - delayed entries: 2000-2013
  - end of follow-up: 2014
- Baseline covariates: age, sex, cardiovascular history
- Longitudinal biomarker (yearly): serum creatinine



# Descriptive statistics & censoring issue

- $s \in \mathcal{S} = \{0, 0.5, \dots, 5\}$
- $t = 5$  years



# Joint model

## ► Longitudinal

$$\begin{aligned} \log \left[ Y_i(t_{ij}) \right] &= (\beta_0 + b_{0i}) + \beta_{0,\text{age}} \mathbf{AGE}_i + \beta_{0,\text{sex}} \mathbf{SEX}_i \\ &\quad + (\beta_1 + b_{1i} + \beta_{1,\text{age}} \mathbf{AGE}_i) \times t_{ij} + \epsilon_{ij} \\ &= \mathbf{m}_i(t) + \epsilon_{ij} \end{aligned}$$

(fitted using  package **JM**)



# Joint model

## ► Longitudinal

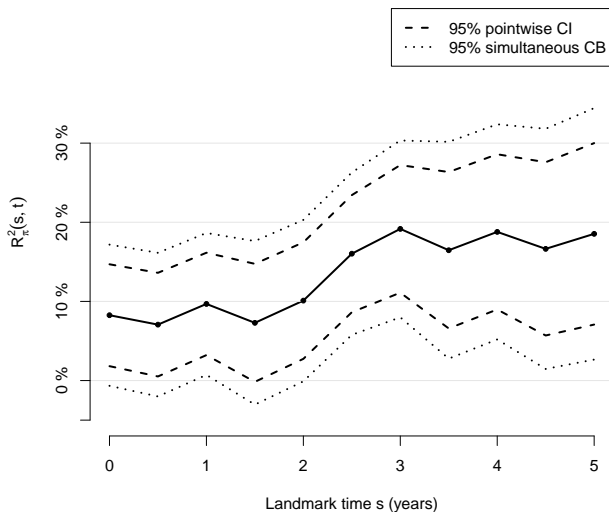
$$\begin{aligned} \log \left[ Y_i(t_{ij}) \right] &= (\beta_0 + b_{0i}) + \beta_{0,\text{age}} \mathbf{AGE}_i + \beta_{0,\text{sex}} \mathbf{SEX}_i \\ &\quad + (\beta_1 + b_{1i} + \beta_{1,\text{age}} \mathbf{AGE}_i) \times t_{ij} + \epsilon_{ij} \\ &= \mathbf{m}_i(t) + \epsilon_{ij} \end{aligned}$$

## ► Survival (hazard)

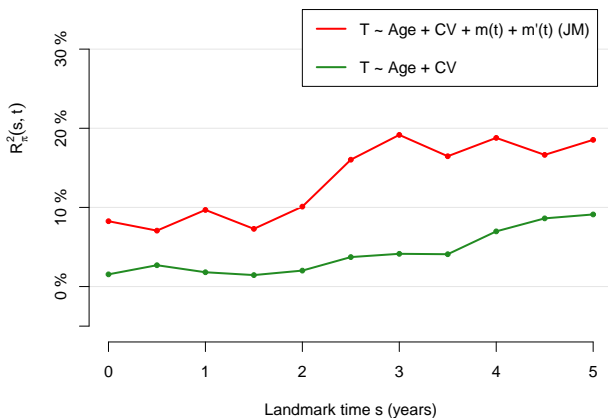
$$\begin{aligned} h_i(t) = h_0(t) \exp \left\{ \gamma_{\text{age}} \mathbf{AGE}_i + \gamma_{\text{cv}} \mathbf{CV}_i \right. \\ \left. + \alpha_1 \mathbf{m}_i(t) + \alpha_2 \frac{d\mathbf{m}_i(t)}{dt} \right\} \end{aligned}$$

(fitted using  package **JM**)



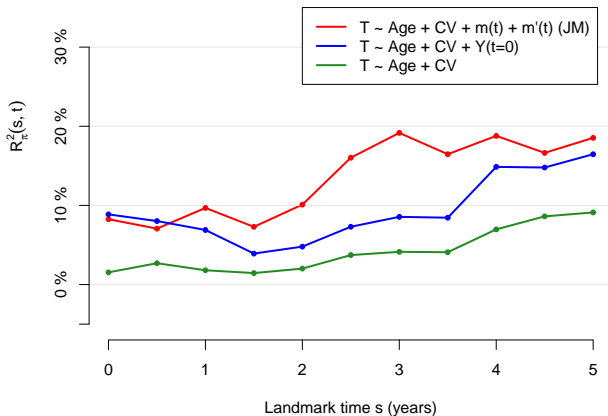
$R_{\pi}^2(s, t)$  VS  $s$  ( $t=5$  years)

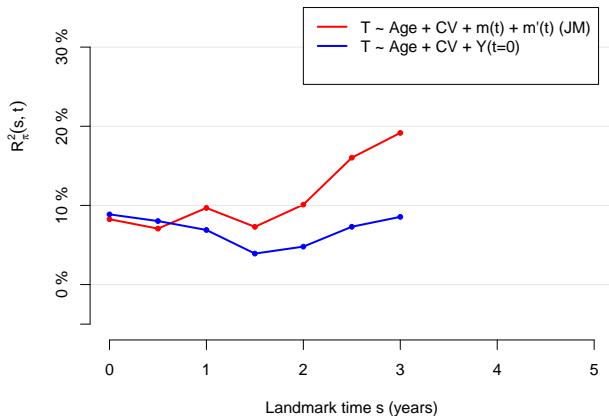
Comparing  $R_{\pi}^2(s, t)$  vs  $s$  for different  $\pi(s, t)$ 

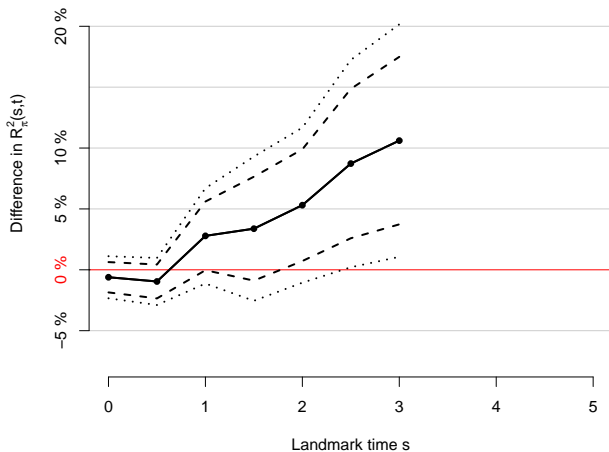
Comparing  $R_{\pi}^2(s, t)$  vs  $s$  for different  $\pi(s, t)$ 



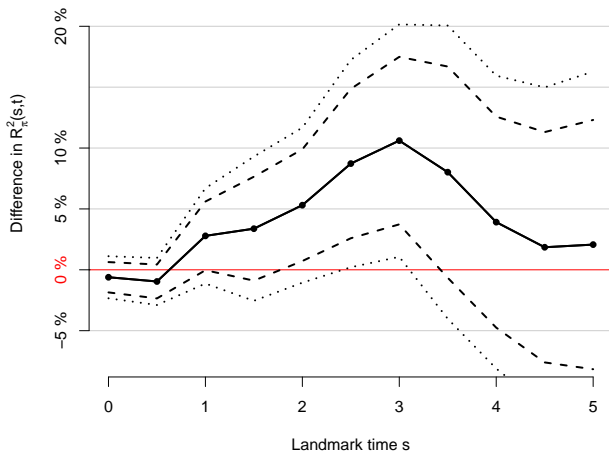
# Comparing $R_{\pi}^2(s, t)$ vs $s$ for different $\pi(s, t)$



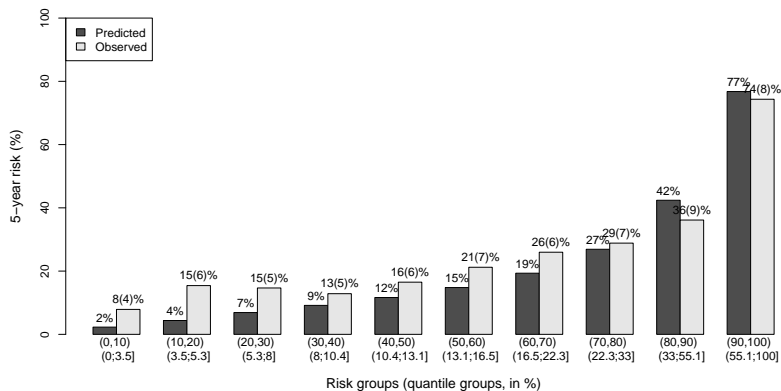
Comparing  $R_{\pi}^2(s, t)$  vs  $s$  for different  $\pi(s, t)$ 

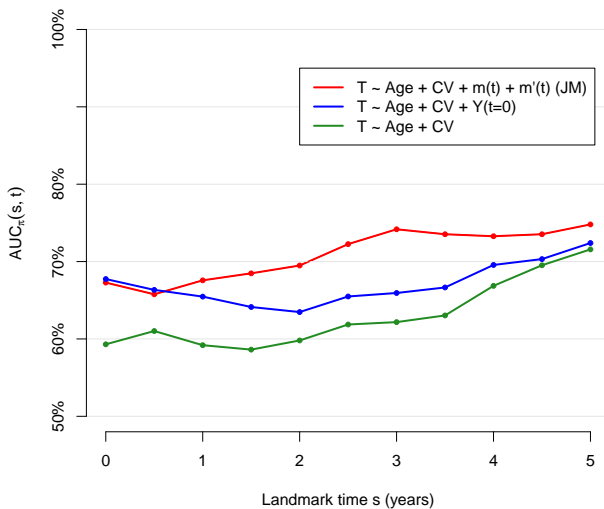
Comparing  $R_{\pi}^2(s, t)$  vs  $s$  for different  $\pi(s, t)$ 

# Comparing $R_{\pi}^2(s, t)$ vs $s$ for different $\pi(s, t)$



# Calibration plot (example for $s = 3$ years)



Area under the ROC( $s, t$ ) curve vs  $s$ 

# Summing up

## ▶ **R<sup>2</sup>-type curve**

- summarizes calibration and discriminating simultaneously
- has an understandable trend

## ▶ **Simple model free inference**

- predictions can be obtained from any model
- we do not assume any model to hold
- allows fair comparisons of different predictions

## ▶ **The method accounts for:**

- Censoring
- Dynamic setting (the at risk population changes)



## Discussion

- ▶ The strong calibration assumption allows different interesting interpretations:
  - Explained variation
  - Correlation
  - Mean risk difference
  
- ▶ Unfortunately
  - the strong calibration cannot be checked (curse of dimensionality)
  
- ▶ However
  - weak and strong definitions are closely related:
    - strong calibration implies weak calibration
    - weak calibration can “often” be seen as a reasonable approximation of strong calibration in practice
  - weak calibration can be assessed (plots)





## Discussion

- ▶ The strong calibration assumption allows different interesting interpretations:
  - Explained variation
  - Correlation
  - Mean risk difference
  
- ▶ Unfortunately
  - the strong calibration cannot be checked (curse of dimensionality)
  
- ▶ However
  - weak and strong definitions are closely related:
    - strong calibration implies weak calibration
    - weak calibration can “often” be seen as a reasonable approximation of strong calibration in practice
  - weak calibration can be assessed (plots)

Thank you for your attention!

