# Reconstruction of Individual Patient Data (IPD) from Published Kaplan-Meier Curves Using Guyot's Algorithm: Step-by-Step Programming in R

Ajay Gupta, Daiichi Sankyo
Natalie Dennis, Daiichi Sankyo

## ABSTRACT

Secondary analysis may require the use of reconstructed patient-level data from published Kaplan-Meier (KM) curves to support a number of different objectives, including indirect treatment comparisons within the context of economic evaluations. Guyot (2012) developed an algorithm that reconstructs individual patient data (IPD) for time-to-event endpoints using published KM curves. This presentation will provide step-by-step instructions and a use case for executing the Guyot (2012) algorithm to reconstruct IPD from published KM curves in R.

## INTRODUCTION

R is a programming language for statistical computing and data visualization. It has been adopted in the fields of data mining, bioinformatics, and data analysis. The core R language is augmented by many extension packages, containing reusable code, documentation, and sample data. R software is open-source and free software. It is licensed by the GNU Project and available under the GNU General Public License. It is written primarily in C, Fortran, and R itself. Precompiled executables are provided for various operating systems. As an interpreted language, R has a native command line interface. Moreover, multiple third-party graphical user interfaces are available, such as RStudio-an integrated development environment.
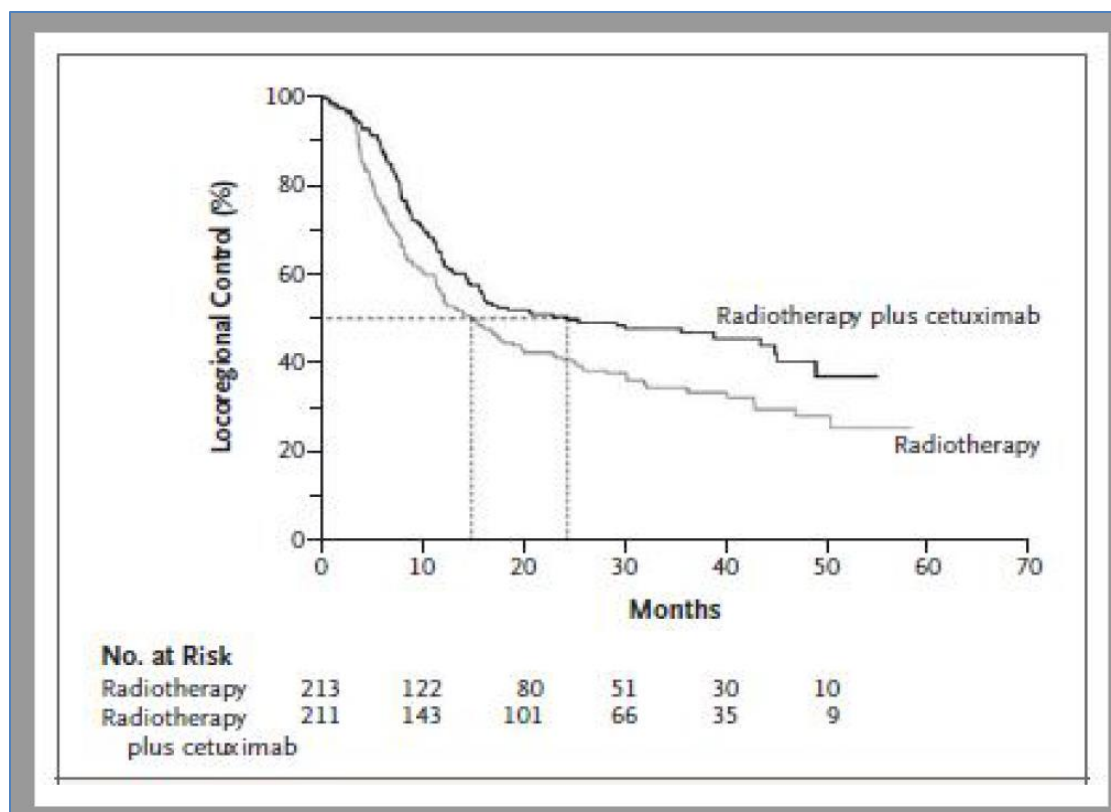
Secondary analysis may require the use of reconstructed patient-level data from published Kaplan-Meier (KM) curves to support a number of different objectives, including indirect treatment comparisons within the context of economic evaluations. Guyot (2012) developed an algorithm that reconstructs individual patient data (IPD) for time-to-event endpoints using published KM curves. This presentation will provide step-by-step instructions and a use case for executing the Guyot (2012) algorithm to reconstruct IPD from published KM curves in R.

## STEPS TO RECONSTRUCT INDIVIDUAL PATIENT DATA (IPD) FROM PUBLISHED KM CURVE

1. Digitize Kaplan-Meier curves using published graph (using PlotDigitizer, GetData Graph Digitizer or other application)

2. Save the extracted survival data (digitized x- and y-coordinates) as a CSV/Excel file.

3. Create a file for the number of patients at risk, including the time points and the lower and upper intervals (based on the digitized Kaplan-Meier curves)

4. Identify the total number of events (if published)

5. Run Guyot's algorithm using R by importing extracted survival data and number of patients at risk files.
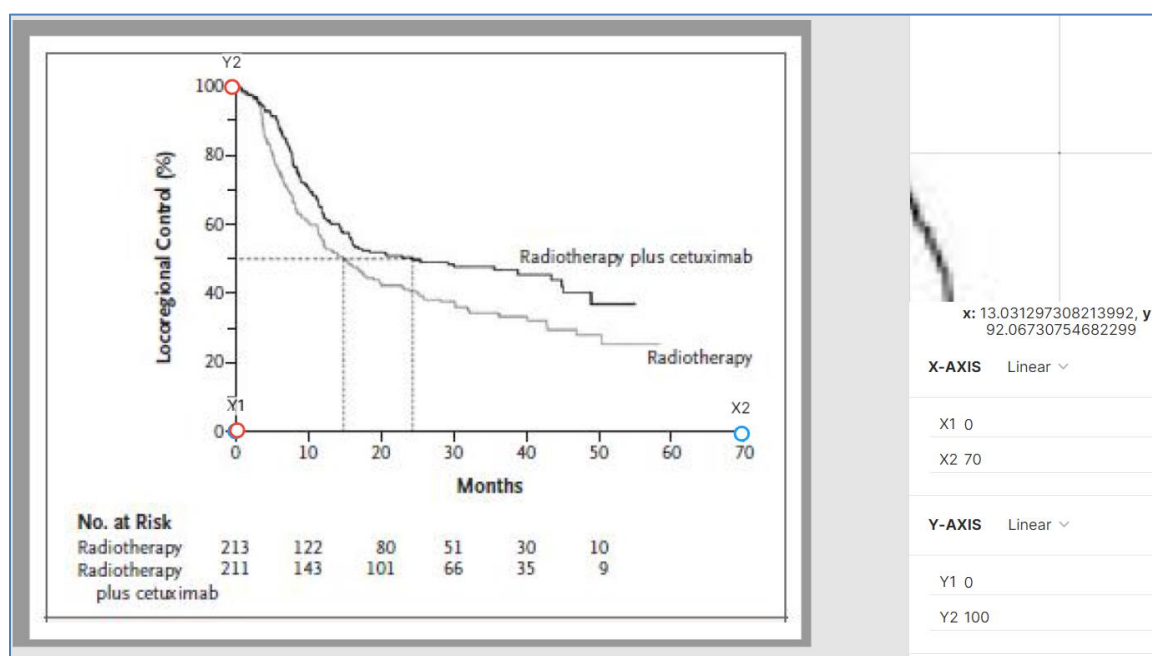
### DIGITIZE KAPLAN-MEIER CURVES USING PUBLISHED GRAPH

1. Create image e.g., PNG, GIF from published KM Curve. See, example from published paper (Guyot 2012).
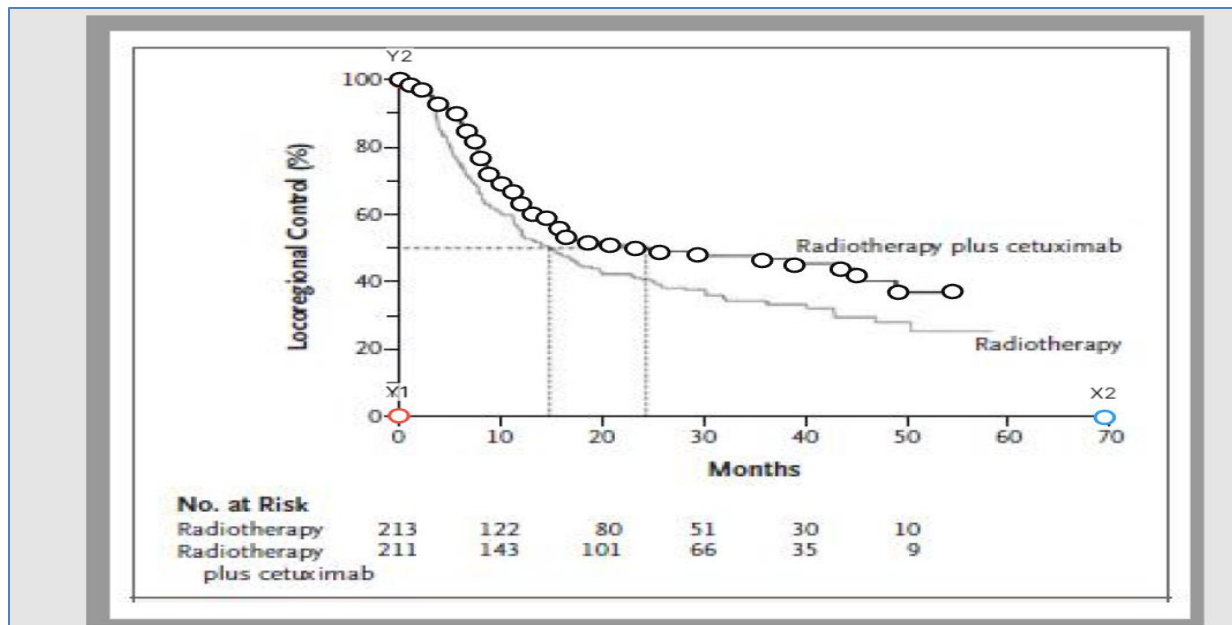
**Display 1. Published KM Curve (Guyot 2012)**

2. Import the file in plot digitizer user interface and select the range for X and Y axis.



**Display 2. Published KM Curve (Guyot 2012) in Plot Digitizer**

3. Mark the event times for the Kaplan-Meier curve to be digitized, i.e., at every point in which there is a step down in the Kaplan-Meier curve.



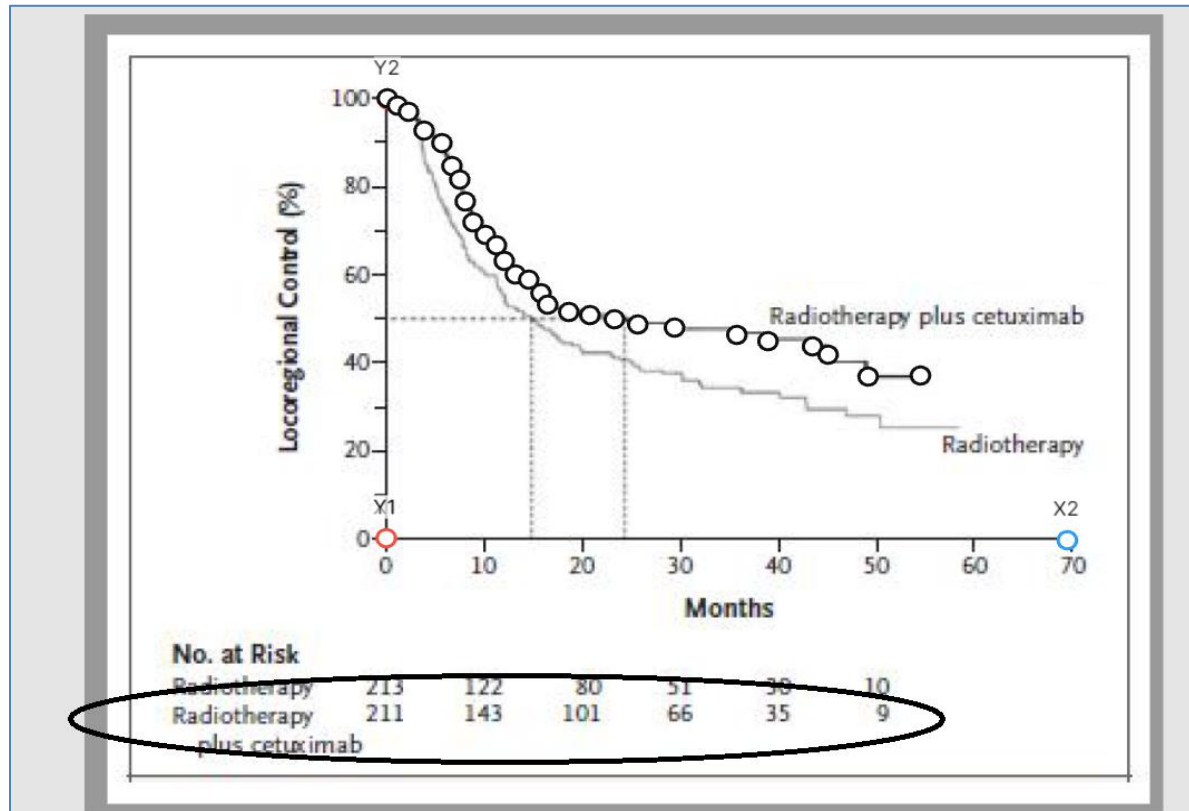**Display 3. Published KM Curve (Guyot 2012) with x and y manual markup.**

4. After marking the relevant coordinates of the Kaplan-Meier curve, export the data into .csv file. Make sure to divide value on Y axis from 100 to get the proportion or mark the Y axis goes from 0 to 1 (even it goes to 100).

| | A | B | C | D |
|---|---|---|---|---|
| 1 | Coordinat | Time | Proportion | |
| 2 | 1 | 0 | 1 | |
| 3 | 2 | 1.143789 | 0.985748 | |
| 4 | 3 | 2.287583 | 0.971496 | |
| 5 | 4 | 3.888887 | 0.928741 | |
| 6 | 5 | 5.718956 | 0.900238 | |
| 7 | 6 | 6.748367 | 0.847981 | |
| 8 | 7 | 7.549024 | 0.817102 | |
| 9 | 8 | 8.120916 | 0.767221 | |
| 10 | 9 | 8.921573 | 0.719715 | |
| 11 | 10 | 10.17974 | 0.691211 | |
| 12 | 11 | 11.32353 | 0.667458 | |
| 13 | 12 | 12.12418 | 0.631829 | |
| 14 | 13 | 13.26797 | 0.60095 | |
| 15 | 14 | 14.64052 | 0.589074 | |
| 16 | 15 | 15.89869 | 0.558195 | |
| 17 | 16 | 16.58496 | 0.532067 | |
| 18 | 17 | 18.75817 | 0.515439 | |
| 19 | 18 | 20.93138 | 0.508314 | |
| 20 | 19 | 23.44771 | 0.498812 | |
| 21 | 20 | 25.84967 | 0.486936 | |
| 22 | 21 | 29.62418 | 0.47981 | |

**Display 4. CSV file from plot digitizer.**

## CREATE A FILE FOR THE NUMBER OF PATIENTS AT RISK:

See below graph and follow the instructions below to create the number of patients at risk file.



**Display 5. Published KM Curve (Guyot 2012) showing number of patients at risk.**

- Create a file for the number of patients at risk, including the time points and the lower and upper intervals.

- Nrisk: value provided in graph above.

- Trisk: time value corresponding to each Nrisk (every 10 months in this example)

- Lower and Upper: the coordinates in the digitize file corresponding to each time window (e.g., coordinates 1-9 fall between 0 and 10 months).

| nrisk | trisk | lower | upper |
|---|---|---|---|
| 211 | 0 | 1 | 9 |
| 143 | 10 | 10 | 17 |
| 101 | 20 | 18 | 21 |
| 66 | 30 | 22 | 23 |
| 35 | 40 | 24 | 26 |
| 9 | 50 | 27 | 27 |

| | A | B | C | D |
|---|---|---|---|---|
| 1 | Coordinate | Time | Proportion | |
| 2 | 1 | 0 | 1 | |
| 3 | 2 | 1.143789 | 0.985748 | |
| 4 | 3 | 2.287583 | 0.971496 | |
| 5 | 4 | 3.888887 | 0.928741 | |
| 6 | 5 | 5.718956 | 0.900238 | |
| 7 | 6 | 6.748367 | 0.847981 | |
| 8 | 7 | 7.549024 | 0.817102 | |
| 9 | 8 | 8.120916 | 0.767221 | |
| 10 | 9 | 8.921573 | 0.719715 | |
| 11 | 10 | 10.17974 | 0.691211 | |
| 12 | 11 | 11.32353 | 0.667458 | |
| 13 | 12 | 12.12418 | 0.631829 | |
| 14 | 13 | 13.26797 | 0.60095 | |
| 15 | 14 | 14.64052 | 0.589074 | |
| 16 | 15 | 15.89869 | 0.558195 | |
| 17 | 16 | 16.58496 | 0.532067 | |
| 18 | 17 | 18.75817 | 0.515439 | |
| 19 | 18 | 20.93138 | 0.508314 | |
| 20 | 19 | 23.44771 | 0.498812 | |
| 21 | 20 | 25.84967 | 0.486936 | |
| 22 | 21 | 29.62408 | 0.47981 | |
| 23 | | | | |

**Display 6. Number of patients at risk file**

## RUN GUYOT'S ALGORITHM USING R:

- Download R program containing Guyot's algorithm from following location.
  12874_2011_700_MOESM1_ESM.PDF (springer.com) and update the R programs with
  respective values (shown in figure below with arrow)

```
#Algorithm to create a raw dataset from DigizeIt readings from a Kaplan-Meier curve

library("MASS")
library("splines")
library("survival")


###FUNCTION INPUTS
path<-"C:\\PHD\\algorithm\\reliability exercice\\"
digisurvfile<-"data initials study2 figA arm1 time1.txt"          #Input survival times from graph reading
nriskfile<-"nrisk study2 figA arm1 time1.txt"                     #Input reported number at risk
KMdatafile<-"KMdata study2 figA arm1 time1 ne.txt"                #Output file events and cens
KMdataIPDfile<-"KMdataIPD study2 figA arm1 time1 ne.txt"          #Output file for IPD
tot.events<-"NA"                          #tot.events = total no. of events reported. If not reported, then tot.events="NA"
arm.id<-1 #arm indicator
###END FUNCTION INPUTS

#Read in survival times read by digizeit
digizeit<- read.table(paste(path,digisurvfile,sep=""),header=TRUE)
t.S<-digizeit[,2]
S<-digizeit[,3]

#Read in published numbers at risk, n.risk, at time, t.risk, lower and upper
# indexes for time interval
pub.risk<-read.table(paste(path,nriskfile,sep=""),header=TRUE)
t.risk<-pub.risk[,2]
lower<-pub.risk[,3]
upper<-pub.risk[,4]
n.risk<-pub.risk[,5]
n.int<-length(n.risk)
n.t<- upper[n.int]
```

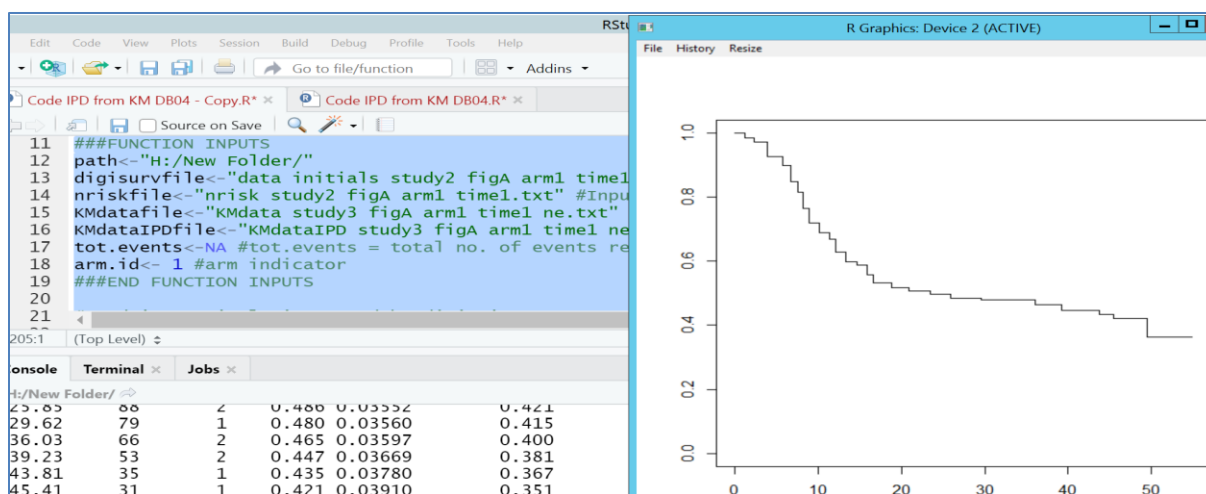**Display 7. R program from the Guyot (2012) publication**

- Execute the R code in R studio:

```
11   ###FUNCTION INPUTS
12   path<-"H:/New Folder/"
13   digisurvfile<-"data initials study2 figA arm1 time1.txt" #Input survival times from graph readi
14   nriskfile<-"nrisk study2 figA arm1 time1.txt" #Input reported number at risk
15   KMdatafile<-"KMdata study3 figA arm1 time1 ne.txt" #Output file events and cens
16   KMdataIPDfile<-"KMdataIPD study3 figA arm1 time1 ne.txt" #Output file for IPD
17   tot.events<-NA #tot.events = total no. of events reported. If not reported, then tot.events="NA
18   arm.id<- 1 #arm indicator
19   ###END FUNCTION INPUTS
20
21   #Read in survival times read by digizeit
22   surv_times <- read.csv("Test_plot_1.csv")
23   digizeit<- data.matrix(surv_times)
24   digizeit[1,2]=0
25   t.S<-digizeit[,2]
26   S<-digizeit[,3]
27
28   #Read in published numbers at risk, n.risk, at time, t.risk, lower and upper
29   # indexes for time interval
30   nrisk trisk <- read.excel("Test nrisk trisk.xlsx")
```

**Display 8. R program in R studio**

- After execution of the programs, to ensure similarity and that the function accurately created pseudo-IPD based on the digitized Kaplan-Meier curve it is possible to recreate the KM curve to compare to the published one.
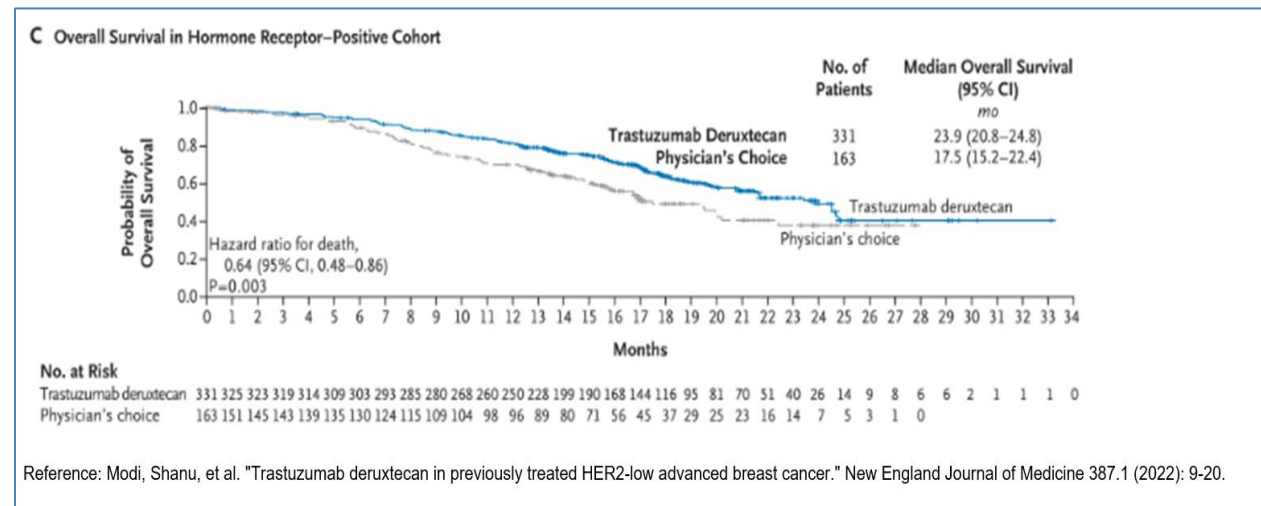


**Display 9. Recreating the KM curve in R studio**

- The programs also create a .txt file with pseudo-patient level data that can be use in secondary analysis.



**Display 10. Individual Patient level data**

## MORE EXAMPLES

See below few more examples. Follow the steps above to recreate the IPD from KM curve.



**C** Overall Survival in Hormone Receptor–Positive Cohort

|  | No. of Patients | Median Overall Survival (95% CI) mo |
|---|---|---|
| Trastuzumab Deruxtecan | 331 | 23.9 (20.8–24.8) |
| Physician's Choice | 163 | 17.5 (15.2–22.4) |

Hazard ratio for death, 0.64 (95% CI, 0.48–0.86) P=0.003

No. at Risk
Trastuzumab deruxtecan 331 325 323 319 314 309 303 293 285 280 268 260 250 228 199 190 168 144 116 95 81 70 51 40 26 14 9 8 6 6 2 1 1 1 0
Physician's choice 163 151 145 143 139 135 130 124 115 109 104 98 96 89 80 71 56 45 37 29 25 23 16 14 7 5 3 1 0

Reference: Modi, Shanu, et al. "Trastuzumab deruxtecan in previously treated HER2-low advanced breast cancer." New England Journal of Medicine 387.1 (2022): 9-20.

**Display 11. Published KM curve.**



| | A Coordinat | B Time | C Proportion |
|---|---|---|---|
| 1 | Coordinat | Time | Proportion |
| 2 | 1 | 0 | 1 |
| 3 | 2 | 0.595812 | 0.993397 |
| 4 | 3 | 1.167244 | 0.993397 |
| 5 | 4 | 2.25816 | 0.993397 |
| 6 | 5 | 2.959463 | 0.981056 |
| 7 | 6 | 3.712714 | 0.972829 |
| 8 | 7 | 4.439991 | 0.972829 |
| 9 | 8 | 4.621811 | 0.960489 |
| 10 | 9 | 5.11532 | 0.960489 |
| 11 | 10 | 5.60883 | 0.952261 |
| 12 | 11 | 5.7387 | 0.948148 |
| 13 | 12 | 6.414028 | 0.944034 |
| 14 | 13 | 6.595848 | 0.931694 |
| 15 | 14 | 6.621823 | 0.923467 |
| 16 | 15 | 7.245203 | 0.911126 |

| | A nrisk | B trisk | C lower | D upper |
|---|---|---|---|---|
| | 331 | 0 | 1 | 3 |
| | 323 | 2 | 4 | 6 |
| | 314 | 4 | 7 | 11 |
| | 303 | 6 | 12 | 18 |
| | 285 | 8 | 19 | 24 |
| | 268 | 10 | 25 | 30 |
| | 250 | 12 | 31 | 35 |
| | 199 | 14 | 36 | 40 |
| | 168 | 16 | 41 | 46 |
| | 116 | 18 | 47 | 53 |
| | 81 | 20 | 54 | 57 |
| | 51 | 22 | 58 | 59 |
| | 26 | 24 | 60 | 62 |
| | 9 | 26 | 63 | 63 |

**Display 12. Markup file and Risk file for above KM curve**

**Display 13. Recreate KM curve using Guyot (2012) algorithm.**



**Display 14. Individual Patient level data**

## ALTERNATIVE APPROACH

- In 2021, Na Liu, Yanhong Zhou & J. Jack Lee proposed a modified, more flexible version of Guyot's algorithm to reconstruct IPD from published K-M curves and developed a R package and Shiny application. See below publication link for more detail. More details will be provided in future presentation.

  - IPDfromKM: reconstruct individual patient data from published Kaplan-Meier survival curves | BMC Medical Research Methodology | Full Text (biomedcentral.com)

- The alternative approach does not require a manual approach to digitizing curves using external software and relaxes some of the requirements for data input. It is therefore an all-in-one approach that reconstructs IPD directly from the KM curve.

## CONCLUSION

Using the digitization software and Guyot (2012) algorithm we can efficiently reconstructs individual patient data (IPD) for time-to-event endpoints using published KM curves. This data can be especially useful in secondary analysis to support a number of different objectives, including indirect treatment comparisons within the context of economic evaluations.

## REFERENCES

R: The R Project for Statistical Computing (r-project.org)

Enhanced secondary analysis of survival data: reconstructing the data from published Kaplan-Meier survival curves | BMC Medical Research Methodology | Full Text (biomedcentral.com)

PNS210 A Comparison of Graph Digitization Software for the Reconstruction of Published Kaplan Meier Curves - Value in Health (valueinhealthjournal.com)

PlotDigitizer Online App

IPDfromKM: reconstruct individual patient data from published Kaplan-Meier survival curves | BMC Medical Research Methodology | Full Text (biomedcentral.com)

R (programming language) - Wikipedia

## ACKNOWLEDGMENTS

## CONTACT INFORMATION

Your comments and questions are valued and encouraged. Contact the author at:

Ajay Gupta, M.S.
Daiichi Sankyo, Inc
211 Mt Airy Rd
Basking Ridge, NJ 07920
Phone : (269) -873-1145
E-mail : Ajgupta@dsi.com,
        ajaykailasgupta@gmail.com

Natalie Dennis
Daiichi Sankyo, Inc
France
Phone : +33 6 77 50 02 77
E-mail : natalie.dennis@daiichi-sankyo.eu

## DISCLAIMER

The content of this paper are the works of the authors and do not necessarily represent the opinions, recommendations, or practices of Daiichi Sankyo.

Any brand and product names are trademarks of their respective companies.