# Help File for KMshiny Application

Na Liu, J. Jack Lee

July 13, 2018

## 1 Description

This package is an implementation to reconstruct individual patients records (time, status) from coordinates extracted from an survival curves. Digitized softwares such as DigitizeIt(for MAC or windows) OR ScanIt(for windows) can be used to get the coordinates from published curves. Other available information like [1] reported patient numbers at risk reported at the beginning of time intervals, total number of initial patients, and total number of events, can also be involved during calculations to increase accuracy.

## 2 Extract Coordinates from published KM curves

### 2.1 Prepare your image

Obtain an image of your graph from PDF, or from your scanner. It does not matter weather the image is capture straightly. However, better resolutions will present better accuracy in the later steps. I will use the published locoregional control events in head and neck cancer KM curve as example.
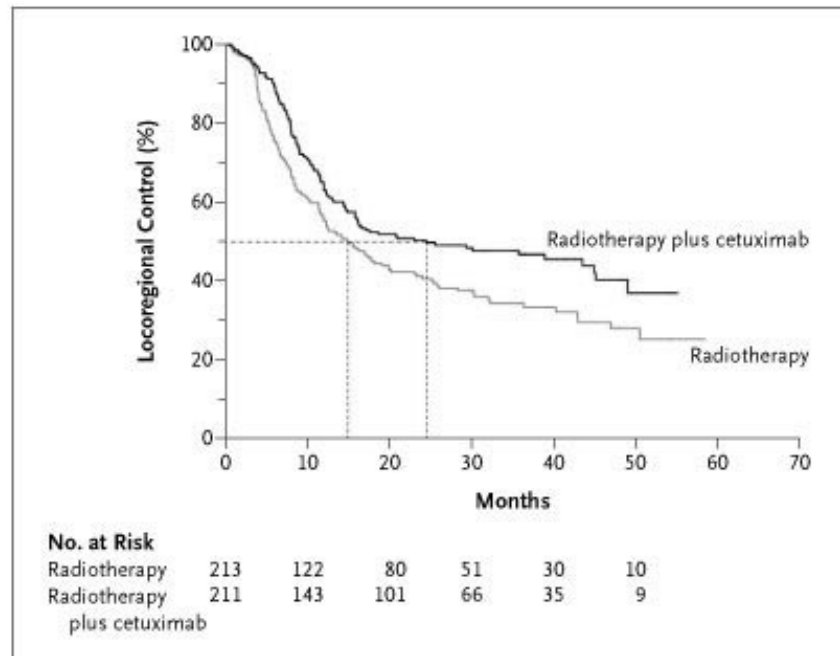
Figure 1: Locoregional control example of published Kaplan-Meier curves. We use the same sample as Guyot's paper to make comparison.

## 2.2 Use ScanIt software to extracting the coordinates of the survival curve.

ScanIt is a free download software to extracting data from scientific graphs, particularly from line and scatter plots.

Once open your file, first we need to identify the axes. We pick two points on each axis. The zoom view will help to pick accurately. Commonly three points are picked: the intersection of the two axes and another point at each axis. Two separate points on each axis are also supported. ScanIt supports linear and logarithmic axes.
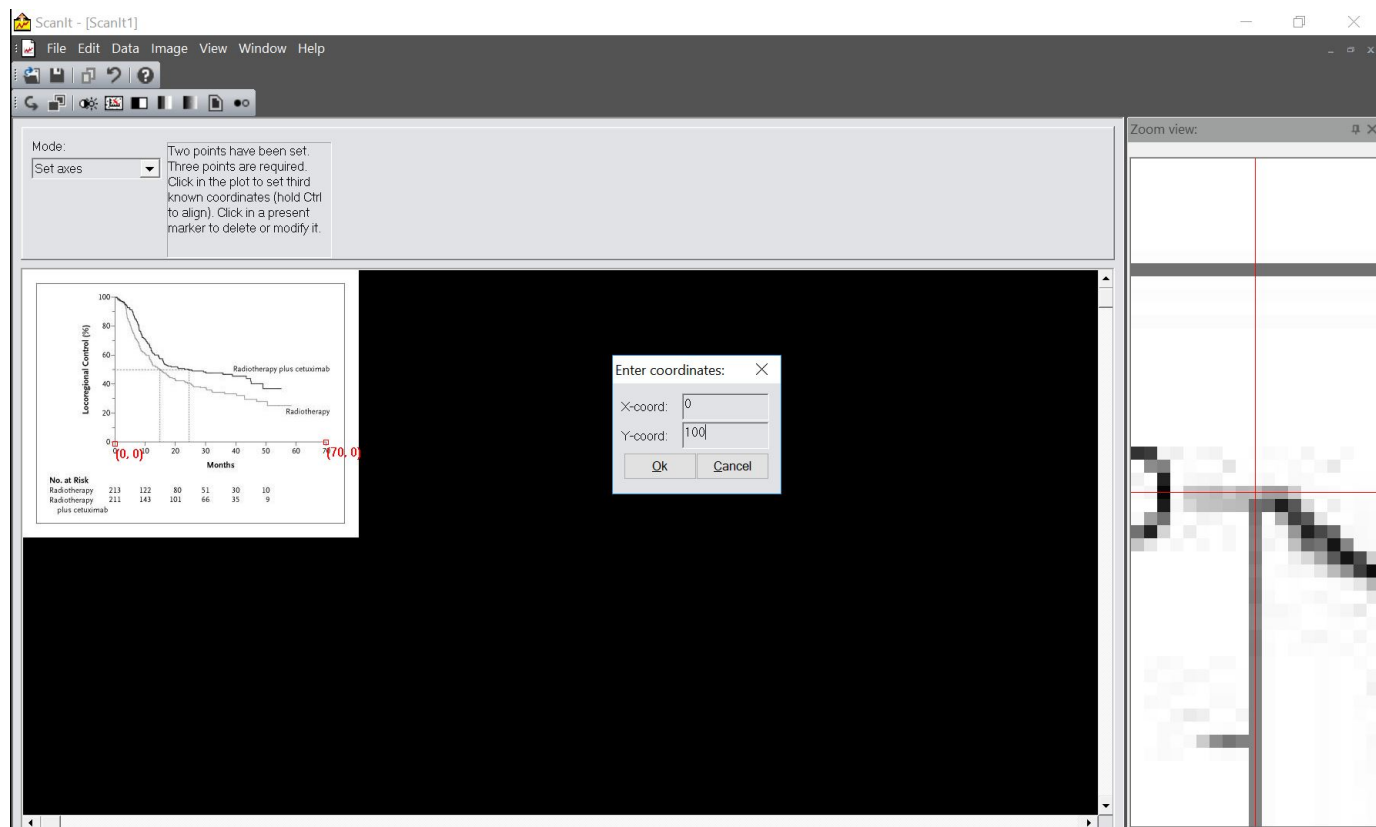
Figure 2: Setting up axis positions
We use zoom view to help pick accuracy.

Then we choose "Trace curve" mode, where entire curves can be digitized by selecting only a few points (2 or more) on the curve. Since the KM curves. generally have a lot of jumps, I normally pick about 10 to 20 points. Again, the zoom view will help the accuracy.
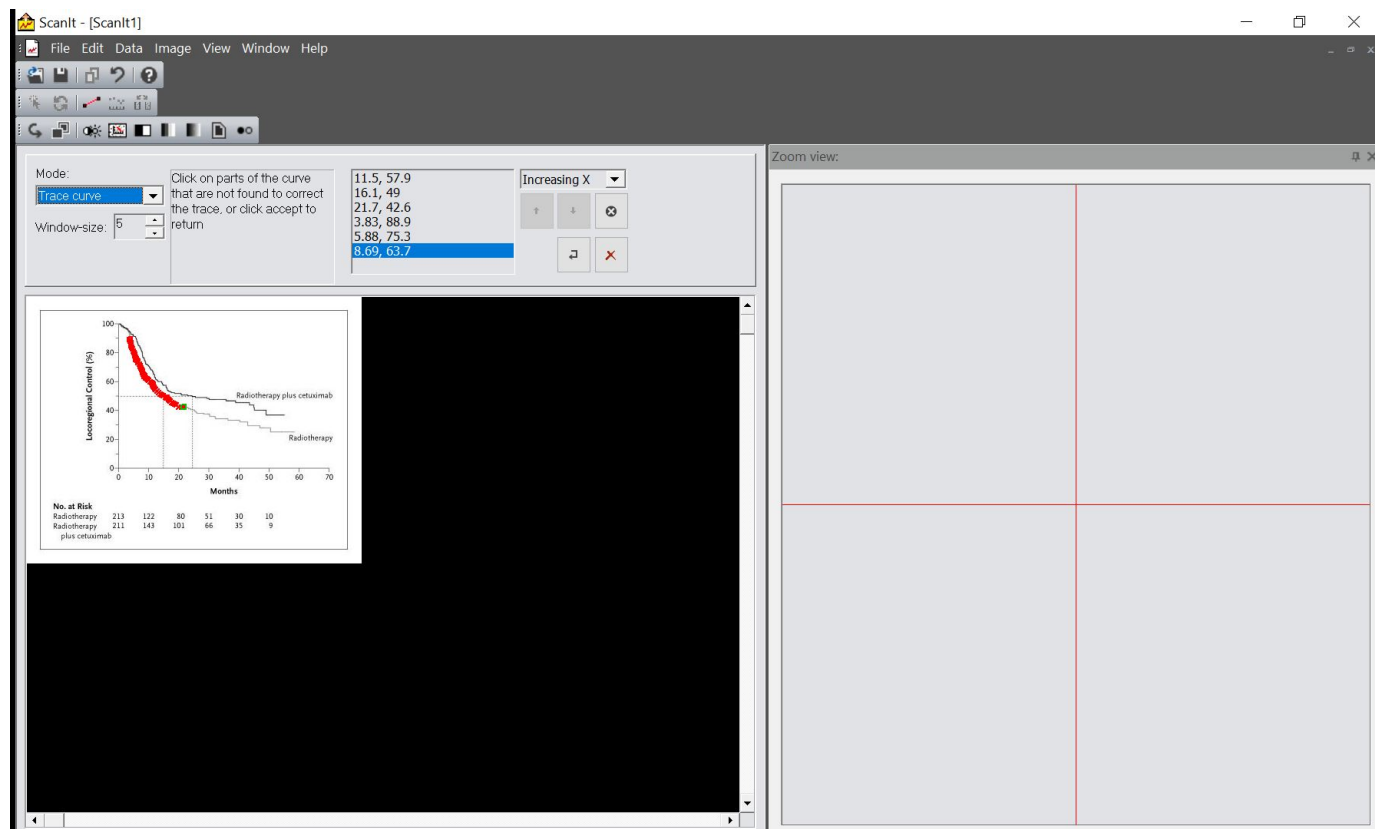
Figure 3: Using Trace curve mode to extract the entire curves automatically
We need to pick about 10 points to tell software the jumps positions.

Then we are read to export the data. The data can be exported as .txt file.
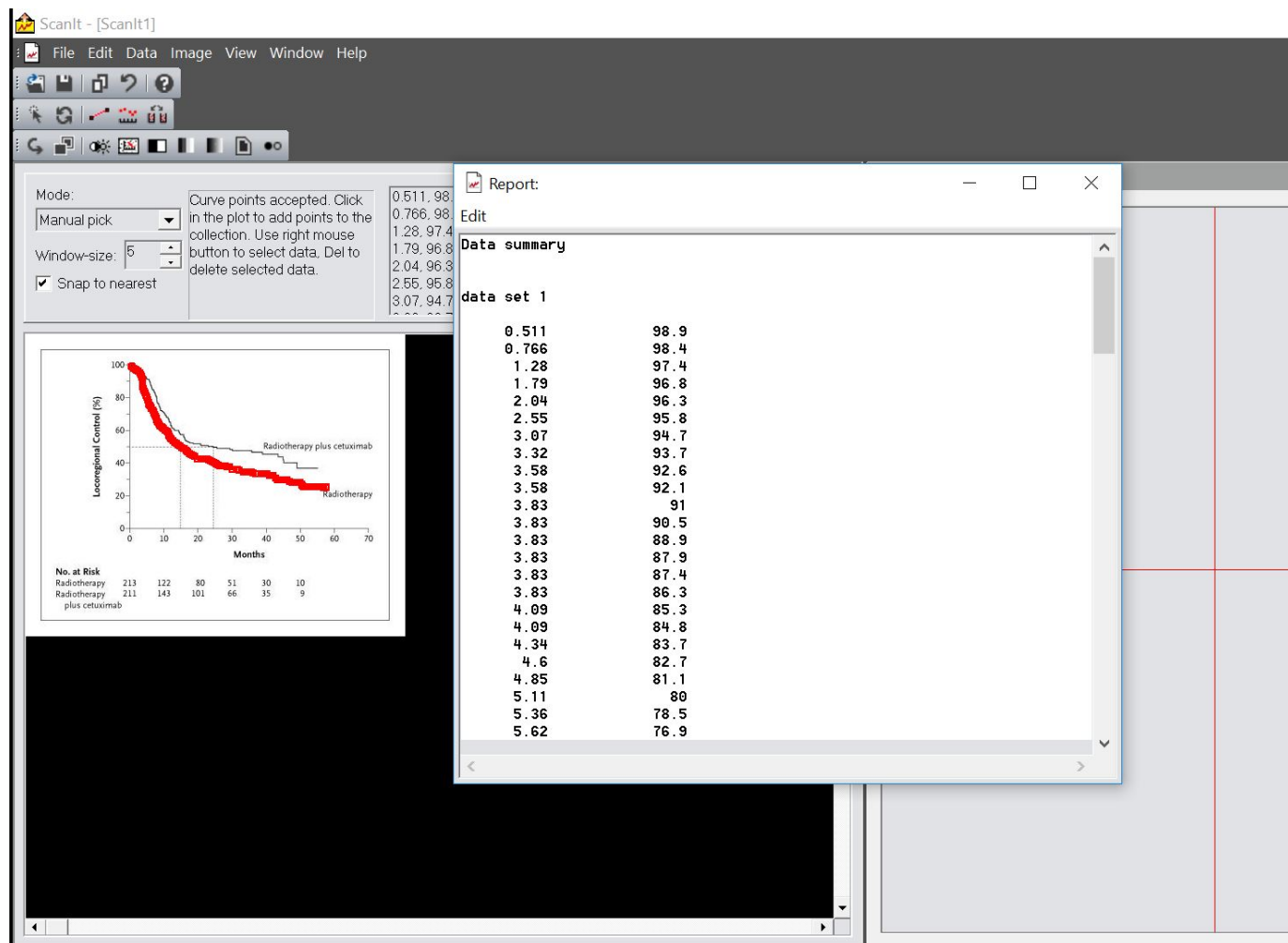
Figure 4: Data can be exported as .txt file

# 3 Using the shiny application to find the individual patients records

## 3.1 When the patient number at risk and time interval known

Most of published KM survival also reported five to ten time points, and the patient number at risk at those time. These can be found under the x-axis of the survival curve.

For example, the radio therapy example we used here, reported at time

$$(0, 10, 20, 30, 40, 50, 60, 70)$$

months, the patient number at risk observed were

$$(213, 122, 80, 51, 30, 10)$$

.

We upload the coordinate datafile we got from ScanIt software, and then input the time at risk and patient number at risk as shown in Figure 5. Other input includs:

- Scale of survival rates. For the case the survival rates reported in the published graph showed as the percentages, we set the scale to 100. For the case using the decimal, we set the scale to 1.

- Treatment arm is the just label for the convenience of next step analysis. You can just pick the right label you want, the default is 1.

- The initial patients is just equal to the patient number at risk at time zero. You can leave it as 'NULL', since no additional information involved in this case.

- Some paper also reported total number of events. Input the number can increase the accuracy of estimation.

Once finish input data, click " Begin Calculate". Click "Reset all" to reset all the input as default status.
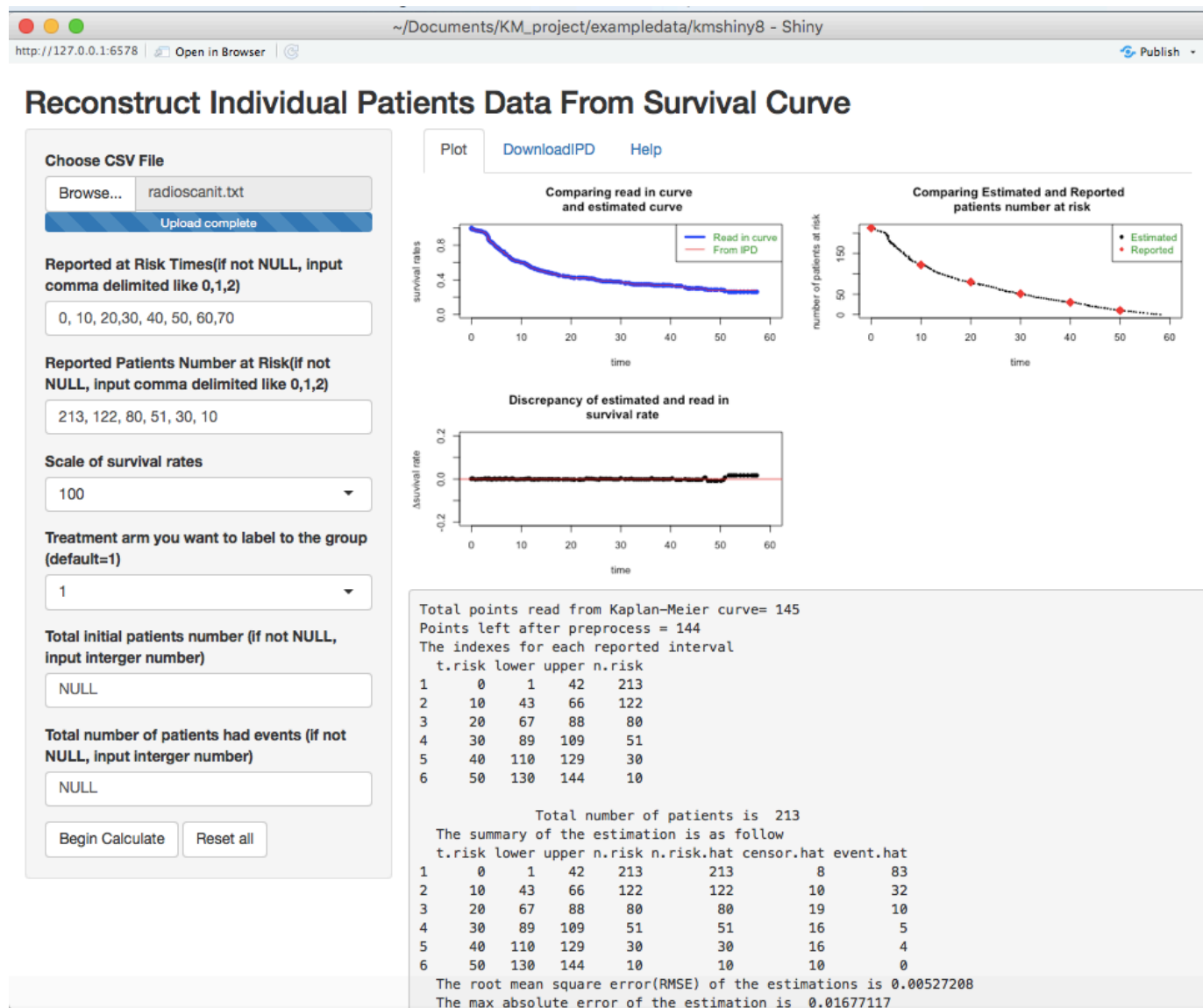
Figure 5: The case when patients number at risk and time interval known

## 3.2 When the patient number at risk and time interval do not known

Occasionally, published survival curve do not reported the patient number at risk and time interval. In this case, you can leave the most of input as the default "NULL" by

click the "Reset all" button. Three data must be input as shown in Figure 6

- the dataset

- Total initial patients number must be input in this case. Or we can not do any estimation.
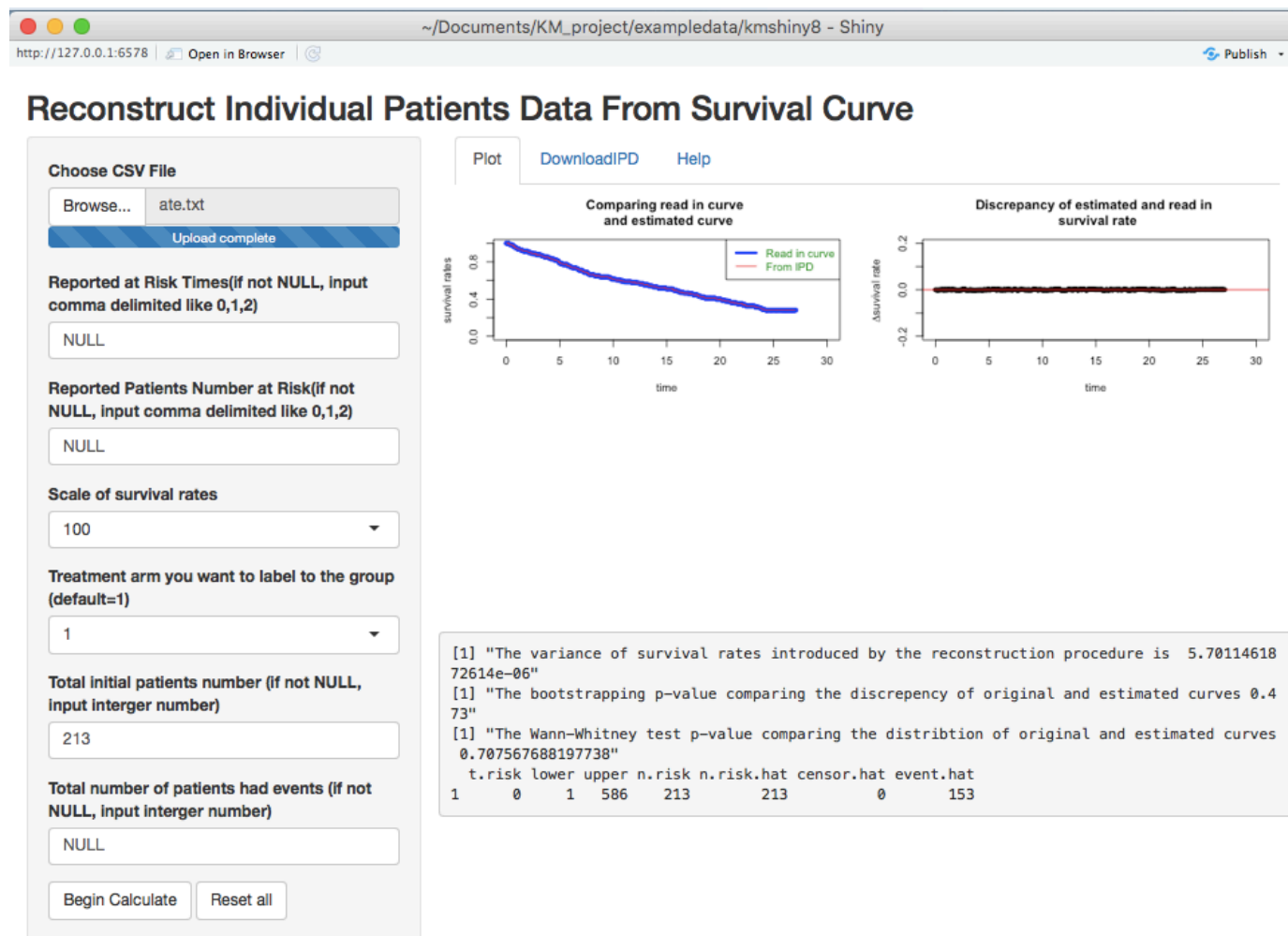
- The scale of survival rates.



Figure 6: The case when patients number at risk and time interval do not know

## 3.3 Explain of the output and figures

We output three figures for the case patient number at risk reported, and two figures for the case patient number at risk did not reported. Those curves shows the accuracy of our estimations.

We performed bootstrapping t-test to comparing the discrepancy of the original read in survival rates and the estimated survival rates. The null hypothesis is the mean discrepancy is zero.

We performed Mann-Whitney test to comparing the distribution of the original and estimated curves. P value reported.

The total variance of survival from our method come from two source. By calculating the mean square error, we have the estimation of the sample variance from the reconstruction. And for the read in precision, we generally read in more than 500 points for each curve. Then the standard error of read in is about

$$SE_{readin} = \frac{1}{500}$$

. So the total variance of survival by the reconstruction procedure is

$$Variance = MSE_{survival} + (\frac{1}{500})^2$$

## 3.4 Download the reconstructed patient record

Click "DownloadIPD" button, the reconstructed patient record will be shown in a table. Each row indicates one patients record, include time, status ( have event label as 1, be censored label as 0), and the treatment arm. Download the .csv file for further study.

# References

[1] Patricia Guyot, A. E. Ades, Mario J.N.M. Ouwens, and Nicky J. Welton. Enhanced secondary analysis of survival data: Reconstructing the data from published kaplan-meier survival curves. *BMC Medical Research Methodology*, 12, 2012.