

Independent Study of Robust estimation in accelerated failure time models

Korede Adegboye

400014008

STATS 756 - Topics in Biostatistics

December 3, 2021

Introduction

When considering lifetime data with covariates, information (characteristics) is given by these covariates that have an affect on an individuals lifetime. When the interest of a subject is survival time, characteristics may include age, sex, severity of disease, smoking status, results of blood tests and other laboratory data. In all, when covariates are introduced to that of a statistical model, it is expected that better production of lifetimes and thus, efficient analysis of the study can take place. Analogous to other statistical models, covariates can be quantitative or categorical. In most cases, the goal is to introduce covariates, to parametric models. Where through analysis, we can verify that the assumed distribution(s) are reliable for the set of covariates given.

Incorporating covariates into the traditional linear regression model, is one of several ways. Here we have that the residuals are normally distributed with mean zero and variance σ^2 , and that the response variable is normally distributed with mean μ and variance σ^2 . When considering the Weibull distribution, there are 2 forms in which a regression model can be formed.

Start by noting the the pdf of the Weibull distribution as

$$S(t) = e^{-(\frac{t}{\alpha})^\eta}$$

The first form and second forms derived from the Weibull distribution are denoted by the following, respectively,

$$S(t) = e^{-(\frac{t}{\alpha(x)})^\eta}$$

Most of the robust developments and the one popular for analysis among those in biomedical sciences make use of the semi-parametric Cox's proportional hazards model. Why ??? In eq 1, the perspective used resulted in proportional hazard, now of consideration is that the defining property is proportional hazards. This means..This has an affect on how the covariates are also introduced. Thus, this gives a regression model for lnh

Note other considerations stem from introducing the covariates such that it acts on different characteristics of a model can, produces a different model. Instead of the proportional hazards function, use proportional odds model to get a linear model for log odds. In other words, the result is logistic regression.

$$S(t; x) = e^{-(te^{\beta'x})^\eta}$$

The effect of the covariates, x is to accelerate or decelerate the time scale t as $(te^{\beta'x})$. For example...

Show a demonstration of how the accelerated life time data can be interpreted...

When the goal of the study is to produce an analysis based on regression models for lifetime data with covariates we can apply the accelerated life time model. The accelerated failure time model describes a situation where the subjects history of an event is accelerated. This subject can be of a biological or mechanical nature. For example,...

Popular in practice because it supports various survival models. Its usually chosen in contrast to proportional hazards models where a distribution may be useful to the model. Its use requires the assumption of a probability distribution. To assume such an application, one must assess which distributions fit the data best. With that it requires precise specifications, otherwise it is expected for it to be an insufficient model. With covariates being the main component of accelerated lifetime models, accelerated lifetime regresses the logarithm of these covariates. Where many find it useful when analyzing censored survival data.

The motivation of this paper, stems from the sensitivity of outliers when modelling. Like traditional regression outliers are of importance to survival analysis. They tend to be more sensitive to the most influential observations. (Pinto JD, Carvalho AM, Vinga S 2015) have researched the detection of outliers in Cox proportional hazard models based on the concordance c-index and named the method Dual Bootstrap Hypothesis Testing (DBHT). The c-index is analogous to the receiver operating characteristic (ROC) curve, where interest is in the sensitivity and specificity of a model. When outliers are present the c-index lowers, and thus suggests the model is sensitive.

Dual Bootstrap Hypothesis Testing is an improvement on Bootstrap Hypothesis Testing (BHT). BHT is the a method that removes a single observation from the data set, and analyzes the c-index to signal outliers

(and inliers). It has been discovered that BHT increases the number of outliers, so DBHT aims to provide a solution. Note that, bootstrapping is a resampling technique that aims to uncover the true distribution given by the data (rephrase a little). DBHT takes two bootstrap samples, where the ‘poison’ sample comes from the original data set and the ‘antidote’ is taken where one observation is removed from the original data set. Both samples are then compared via p-values. Essentially this is a test to see if an observation is “poison” (outlier). The test between the two samples gives a null hypothesis that the expected c-index of “antidote” sample is greater than that of the “poison” sample. Whereas the alternative hypothesis would be that expected c-index of the “antidote” sample is less than or equal to that of the “poison” sample. The p-values are calculated via the two sample t-test (Welch’s t-test) for unequal variances.

where unlikely observations are rejected or downweighted

Show some graph...

The main problem that is being raised in this study is that outliers in cox regression models result in sensitive estimators. Note outliers tend to tell truth about the data, and it may be necessary to investigate what may have happened when the data was recorded. Nevertheless, investigations could be lengthy or inconclusive so proceeding with robust estimation via robust maximum likelihood estimator should be feasible so long as its use gives a model where outliers have been detected and are given less weight. In the past 30 years, there has been extensive development regarding robust methods that utilized semi-parametric cox proportional hazards model. Generally speaking, these methods did not effectively handle covariate outliers. Dr. Sinha proposes the use of parametric accelerated failure time model. With the aim to produce a parametric accelerated failure time model that is able to handle outliers in 2 settings – survival outcome variable and covariates.

Robust against outliers in survival times but not against potential outliers in covariates

Model, notation, and method

The log-linear form of the accelerated failure time model

Accelerated failure time models include exponential, Weibull, log-normal, and log-logistic distributions. - The purpose for using an accelerated life time model stems from the assumption that the covariates on a subject are accelerated (tends to multiply) in regards to time rather than hazards (proportional hazards model). In others accelerated lifetimes can be applied when interest in analyzing the speed of progression of a disease over time. Equation () above gives the general accelerated failure time model. To derive the log-linear form,

consider the parametric accelerated failure time model. It can be denoted by the following...

Graphs of exponential,

Interest lies in estimating the model parameters given by the parametric distribution in such a fashion that it is resistant to outlier. - Show formulas...

Robust Estimation

Robust estimation can be defined as a method that is not influenced by slight differences of the proposed assumptions (ie., probability distribution). As a consequence it optimizes the algorithm's search for a robust estimator. In other words, the results obtained are deemed to be reliable given these slight variations. In general, maximum likelihood estimates, linear combinations and statistical rank tests can aid in achieving robust estimation. Here, of interest is estimating the

Estimations of parameters can be derived via the standard maximum likelihood method. The i -th δ_i denotes the censoring status of the i -th survival time. From here, inference can proceed as follows from the maximum likelihood theory. It is the likelihood score that is a function of error terms and covariates that are unbounded for both the survival time t_i and covariates x_i . Here it has been determined (by which result???) that this construction is not robust to outliers (sensitive to outliers). Therefore the motivation of this literature comes from this issue. Dr. Sinha proposes that the solution for this includes the Huber-type monotone psi function on the error terms to bound the influence of outliers. - Explain more about Huber-type monotone psi function, (Huber 1973)

Here our goal is to estimate the model parameters, and by a suitable robust method which is resistant to potential outliers in both survival times and covariates in the data

Asymptotics

- Sketch of asymptotic properties. Recall, by the term asymptotic leads to determine the properties that are assumed to correct when dealing with a large sample. Expectation using the mean value theorem. Normality if the estimator is properly behaved, then the central limit properly can be applied. Covariance

Robust estimation in Weibull accelerated failure time models

- Review of the log-linear form of the accelerated failure time model. With the baseline hazard function in relation to the Weibull distribution the scale and shape parameters are constructed differently. More specifically, with log-linear form, we can visualize the error terms as the Gumbel distribution. The Gumbel can be denoted as an extreme value (type 1) and is a member of the exponential family. Traditionally its purpose lies with modelling the distribution of the maximum number of samples of various distributions. In relation to the Weibull, it is the log-weibull distribution. When latent (inferred) variables are considered it has then been discovered by Gumbel that this is in fact the error terms of the log-linear model.

Robust maximum likelihood is used to numerically solve the estimating equations. Recall, that iterative methods include Newton-Raphson, etc. The central idea behind iterative methods is the use of an initial value and subsequent approximations that improve at each iteration until a specified error bound is met. Here other iteration methods could have been considered. For instance, secant method, newtons method, or Bisection method. There are some characteristics that distinguish the Newton-Raphson method from the others. First being that it deals with extrapolation, in which the goal of this survival analysis is to predict values that fall outside the range of data. Which is clear by definition of a accelerated failure time model, where one is interested in analyzing the affects of covariates on a subjects lifetime (ie., how much longer they will live given qualitative or quantative characteristics). Next, the Newton-Raphson method is fast and a standard when the first and second derivatives are convenient to find and compute. Such a method is called a root finding algorithm where of interest is finding a value such that it takes the function to be zero. This is consistent to the motivation behind accelerated lifetime models. (rephrase some parts)

o The initial estimates are Derived from first-order Taylor series expansion is used to derive the initial estimates, and thus the iterative equations until convergence. As a consequence, a $(p+2) \times (p+2)$ matrix is produced for the objective function. It is the equations given by the components, where the huber psi function is imposed. Note that it can tuned. o Approximate variance obtained from the sandwich-type variance-covariance matrix. ie,

$$V(\hat{\theta}) = M^{-1}QM^{-1}$$

where ...

Robust estimation in log-normal accelerated failure time models

- Under log-normal settings there are some differences. Under log-normal settings, the accelerated failure time models have an associated shape and scale parameter. The error terms follow a standard normal distribution. The standard maximum likelihood estimators are displayed, then the robust approach. Which thereafter can be solved numerically using the Newton-Raphson method. Expectations and covariance are in respect to the standard normal distribution
- Generalized M (GM) estimating equations for ordinary linear regression???

Robust estimation in log-logistic accelerated failure time models

- Eq 9 used reference in all considerations of distributions and the equation is solved (in relation to the distribution)
- What does one need to understand in order to grasp the use of RML

Simulation study

Two sets of data are compared against each other. A set in which no outliers were considered in the data. Another set, where the data is altered and polluted with outliers (additional noise). Next of consideration is the performance of these sets under ML and RML. An informal way of detecting outliers is by a residual plot. Both residuals plots clearly show four outliers are present. The plots suggests RML is not as influenced to these outliers as much as that for estimation under ML. The actual values tend to be more than the predicted values (positive trend).

To formally access whether this interpretation is feasible, empirical results such as biases, mean squared errors, coverage probabilities and average lengths of 95% confidence intervals are closely examined using Monte Carlo simulations. Note, coverage probabilities refers to the proportion random regions (intervals) in regards to time that will contain the true value.

First, the empirical results are analysed in Weibull models. When there are no outliers, the maximum likelihood (ML) estimator has a lower mean-squared error than the robust maximum likelihood (RML). This is okay (not a cause for concern), since the focus is on data with outliers – expected in practice – Table 1. In contrast, one can argue that this is a cause for concern, which implies the importance of conducting such a comparison between the two methods. As it should not be overlooked that there exists a number of data that where ML outperforms RML by larger margins. Adding a single outlier clearly suggests that ML begins to produce larger mean squared errors, larger biases, and lower coverage probabilities than the estimators

produced by robust maximum likelihood (RML). Now with 4% (of $n = 100$) outlier, RML has some biases but it is much more severe in the case for ML. When the sample size is increased, it can be shown that the respective estimators do much better under weak law of large numbers (WLLN) but RML generally outperforms the ML method. Recall, WLLN refers to the average of a sequence (large samples) of independent and identically distributed (iid) random variables with a common mean and variance that converges in probability to the true value. - Tables 3 and 4 In regards to log-normal models, it is again determined that when outliers are assumed to present in the data, then the RML outperforms the ML method., - Table 5 and 6 Considering the log-logistic models, the conclusions mentioned for the other two models are consistent. More specifically, the ML method of estimating the first regression parameter is not as sensitive to outliers. When the sample size is increased, although the bias and MSE decrease at a constant rate for both estimators, the RML exceeds the ML. To remark, a reviewer provided feedback that distinguished between moderate and extreme outliers. This insight can be interpreted as the question; how well does the robust method in this paper perform under extreme outliers. Real-life data can contain any number of outliers. More plausible would be a number of moderate outliers and few extreme outliers. Knowing the amount of moderate or extreme outliers could help in the construction of the RML. To address this concern, Dr. Sinha suggests that the robust method would remain more reliable to that of the standard ML.

```
print("this")
```

```
## [1] "this"
```

Application: breast cancer data

In the data that was given for example, there exists ... for quantitative and ...for categorical - If I were to use this data. What else would I need to know. Or can I provide a review on how I used it given the information

Discussion

When choosing between which models how is residual analysis done with applied data - What other determining factors can be used to determine the model? - The RML was compared against ML. But why not compared against the RML proposed by Locatelli et al ? Wouldn't this be the comparison that should be discussed in future ? - I can start thinking myself as to how this robust estimation can be extended to frailty models

References

- Collett D. 2014. *Modelling Survival Data in Medical Research, 3rd Edn.* Chapman; Hall/CRC, New York.
<https://www.routledge.com/Modelling-Survival-Data-in-Medical-Research/Collett/p/book/9781439856789>.
- Huber PJ. 1981. *Robust Statistics*. Lifetime Data Anal 25, 52–78. https://link.springer.com/referenceworkentry/10.1007%2F978-3-642-04898-2_594.
- Lin DY, Wei LJ. 1989. *The Robust Inference for the Cox Proportional Hazards Model*. J Am Stat Assoc 84:1074–1078. <https://www.tandfonline.com/doi/abs/10.1080/01621459.1989.10478874>.
- Pinto JD, Carvalho AM, Vinga S. 2015. *Outlier Detection in Cox Proportional Hazards Models Based on the Concordance c-Index*. Machine learning, optimization,; big data: lecture notes in computer science, pp 252–256. https://link.springer.com/chapter/10.1007/978-3-319-27926-8_22.
- Sinha, S.K. 2019. *Robust Estimation in Accelerated Failure Time Models*. Lifetime Data Anal 25, 52–78. <https://doi.org/10.1007/s10985-018-9421-z>.
- Sinha SK, Rao JNK. 2009. *Robust Small Area Estimation*. Can J Stat 37:381–399. <https://onlinelibrary.wiley.com/doi/10.1002/cjs.10029>.