# Penalized regression for Left-truncated & right-censored survival data

McGough, Incerti, Lyalina, Copping, Narasimhan, Tibshirani (2021) "Statistics in Medicine)

{ A summary by Audrey Yeo Te-Ying).

Rationale: For some inclusion criterions for time to survival data, one sometimes include only patients with a genomic test in a personalised health care context (often the case for oncology medicine that this context is applied). plus with electronic health records (aka, just lots of variables) some high dimensional data (because they are)

(Problem $\longrightarrow$ Solution proposed)

High dimensionality $\longrightarrow$ Penalised Regression such as lasso($l_1$) and ridge ($l_2$), elastic net regression, smoothly clipped absolute deviation

Left truncation at time $t$ $\longrightarrow$ Partial likelihood at time $t$ : $\prod_i^m \frac{e^{x^T B}}{\Sigma e^{x^T B}}$

Overfitting $\longrightarrow$ cross validation

audreytyeo@gmail.com

Things to avoid ——▶ Preventative Measure in
                              Analysis

"immortal time bias" ——▶ ⎛ that's a hard one
because observed patients    ⎜ to avoid when
cannot die prior to          ⎜ "letting" patients
entering the study.          ⎝ into a survival analysis ⎞
                                study                   ⎠

Poor prediction accuracy  ——▶ Don't over fit
limited + generalizability     add regularization
                               penalty   to constrain
                               the size of B,
                               reduce complexity
                               of model

                               avoid p predictors >
                               n sample (p < n
                               is desirable)

audreytyeo@gmail.com

Simulation study (methods)
1) Data generation of } from Weibull.
T_i = latent survival time } } assumes they are uncorrelated
U_i = Right Censoring
V_i = study entry time

observed survival time $= Y_i = \min(T_i, U_i)$

R censor $= T_i > U_i$

L truncated $= V_i > Y_i$

2) Simulation of 10 binary predictors to an $n \times p$ matrix called X. (11 predictors are taken from real-world dataset already, s.t. $p = 21$)

$$\begin{pmatrix} \tilde{X}_{i,1}^* \\ \vdots \\ \tilde{X}_{i,p}^* \end{pmatrix} \sim N\left( \begin{pmatrix} \mu_1 \\ \vdots \\ \tilde{\mu_p} \end{pmatrix}, \Sigma \right) \} \text{ multivariate probit approach}$$

$T_i$, latent survival is taken from Weibull survival function with density

$$f(t) = a m t^{a-1} e^{m t^a}$$

where $m = \exp(\alpha + x^T B)$

s.t. $T_i \sim \text{Weibull}(\alpha, m_i)$

andreytyeo@gmail.com

choice because when using Nelson-Aalen estimator of $\Phi(t)$, cumulative hazard, overall survival is monotonically decreasing, so Weibull can capture baseline hazard

$U_i$, right censoring time is taken also from a Weibull distribution

$V_i$, study entry times is generated in 2 parts, first is the bernoulli random variable $d_i \sim$ Bernoulli$(\tilde{\pi})$ to input if $V$ has status 1 or 0.

$V \sim$ lognormal $(\mu, \sigma^2)$, median = 1 year, mean = 1.6 years

not sure rational for choice
maybe convenience from Normal distributions
as inverse cdf and pdf can be "easily" computed
and (R) skewed (see mean, median chosen)

3) Cross validation with 10 folds where test data is 25% and training is 75% from ~4500 observations of real-world data

3a) repeated 200 times
3b) from small and large $p$ scenarios
$p = 21$     $p = 1011$

4) find $\lambda$ that minimized partial likelihood deviance in (3)

5) plot calibration curves of survival probability of those with and without calibration curves. for small and large $p$ scenario

audreytyeo@gmail.com

6) Use C-index to compare fitted models with and without Left truncation adjustment. the ↑ C-index, the better the model is to discriminate prognosis (survival probability) between patients

7) with penalised Cox proportional models, Compare hazard ratios between models that were and were not adjusted for Left truncation

<span style="color:red">Results (1 example)</span>

C-index

| Model | Covariate size | adjustment | no adjustment |
|---|---|---|---|
| Cox | small | 0.649 / 0.648 not better | 0.580 / 0.580 no diff |
| Cox (lasso) | small | | |
| Cox | large | 0.625 / 0.663 better | 0.556 / 0.600 better |
| Cox (lasso) | large | | |

andreytyco@gmail.com

Points of Discussion

→ generally C-index is higher in model that does not adjust for Left truncation? maybe because of **bias** in test = training sets. therefore careful understanding of the data generation process is crucial ⚠️

→ this approach opens questions for future specification, interpretation and evaluation of prognostic survival models

→ Calibration curves > C-index, maybe latter is more difficult to differentiate between models

→ define risk set appropriately, difficult wh ⓛ truncated data ∴ when and how patients came to be included in the data

→ potential selection bias because only stageIV diagnoses was included eg temporal selection bias, immortal time bias. Can be mitigated by studying the association between left truncated time and survival time.

audreytyeo@gmail.com