

0.1 Regularized AFT models

A standard AFT model is defined as follows:

$$\log(Y_i) = \beta_0 + x_i^T \beta + \sigma \epsilon_i \quad (1)$$

Where x_i are the covariates, Y_i is the observed time (output). $\epsilon_i \sim F$, where F is the distribution function.

Consider F to be the logistic cdf, given as:

$$F(x) = \frac{1}{1 + e^{-x}} \quad (2)$$

The probability density function is given as:

$$f(x) = \frac{e^{-x}}{(1 + e^{-x})^2} \quad (3)$$

Survival function can be calculated from the above:

$$1 - F(x) = \frac{1}{1 + e^x} \quad (4)$$

Here on, we assume that σ is constant for each observation i , and is ignored, and that the covariate matrix is appended with a column of ones. Hence,

$$\epsilon_i = \frac{\log(y_i) - (x_i^T \beta)}{\sigma} \quad (5)$$

For interval regression with censored data, we are given time intervals $\{\underline{t}_i, \bar{t}_i\}$ and covariates x_i for $i = 1 : n$, where \underline{t}_i may be $-\inf$ (left censoring) and \bar{t}_i may be \inf (right censoring).

0.1.1 Likelihood

For calculating likelihood, in the observations with no censoring, the pdf is used, and in censored observations, the cdf is used. The likelihood is given as:

$$L(\beta) = \prod_{i=1}^n \zeta(\beta, \underline{t}_i, \bar{t}_i) \quad (6)$$

where,

$$\zeta_i = \zeta(\beta, \underline{t}_i, \bar{t}_i) = \begin{cases} \zeta_{i,1} = F(R_i) & \text{if: } -\infty = \underline{t}_i, \bar{t}_i < \infty \\ \zeta_{i,2} = 1 - F(L_i) & \text{if: } -\infty < \underline{t}_i, \bar{t}_i = \infty \\ \zeta_{i,3} = F(R_i) - F(L_i) & \text{if: } -\infty < \underline{t}_i \neq \bar{t}_i < \infty \\ \zeta_{i,4} = d(R_i) & \text{if: } -\infty < \underline{t}_i = \bar{t}_i < \infty \end{cases} \quad (7)$$

where, $R_i = \log(\bar{t}_i) - \eta_i$, and $L_i = \log(\underline{t}_i) - \eta_i$. $\eta_i = x_i^T \beta$ is the vector of linear predictors. Clearly, for left censored observations, $\underline{t}_i = -\infty$, hence we can substitute $L_i = 0$ in ζ_3 to get ζ_2 . Similarly, for right censored observations, $\bar{t}_i = \infty$, hence we can substitute $R_i = \infty$ in ζ_3 to get ζ_1 .

Hence the log likelihood is given as:

$$l(\beta) = \sum_{i=1}^n \log \zeta_i = \sum_{i=1}^n \gamma_i \quad (8)$$

where, $\gamma_i = \log \zeta_i$.

0.1.2 Score

For now, we consider that $\underline{t}_i = \bar{t}_i$ is not possible. Hence we are left with the types of censoring.

We consider the partial derivative of $\gamma_{i,3}$ wrt to the parameters β . We can simply obtain derivatives for $\gamma_{i,1}$ and $\gamma_{i,2}$ by setting $L_i = 0$ and $R_i = \infty$ respectively.

We use the following:

$$\frac{\partial F(R_i)}{\partial \beta_j} = -x_{ij} F(R_i) [1 - F(R_i)] \quad (9)$$

$$\frac{\partial \gamma_{i,3}}{\partial \beta_j} = -x_{ij} \frac{[F(R_i)(1 - F(R_i)) - F(L_i)(1 - F(L_i))]}{F(R_i) - F(L_i)} = x_{ij} [F(L_i) + F(R_i) - 1] \quad (10)$$

Define $\mu_i = (F(L_i) + F(R_i) - 1)$, $\mu = [\mu_1, \dots, \mu_n]^T$. Then, partial derivative of log-likelihood is given as:

$$\frac{\partial l(\beta)}{\partial \beta_j} = \sum_{i=1}^n x_{ij} \mu_i \quad (11)$$

Hence, the score (gradient of likelihood) is given as:

$$g = \nabla_{\beta} l(\beta) = X^T \mu = \sum_{i=1}^n \mu_i \bar{x}_i \quad (12)$$

0.1.3 Hessian

Hessian is the second derivative of loglikelihood:

$$H = \frac{\partial g(\beta)^T}{\partial \beta} = \sum_{i=1}^n (\nabla_{\beta} \mu_i) \bar{x}_i^T \quad (13)$$

Using 9, we have:

$$\frac{\partial \mu_i}{\partial \beta_j} = -x_{ij} [F(R_i)(1 - F(R_i)) + F(L_i)(1 - F(L_i))] \quad (14)$$

$$\nabla_{\beta} \mu_i = -\bar{x}_i [F(R_i)(1 - F(R_i)) + F(L_i)(1 - F(L_i))] \quad (15)$$

Hence, the hessian is given as:

$$H = \sum_{i=1}^n -w_i \bar{x}_i \bar{x}_i^T \quad (16)$$

where, $w_i = [F(R_i)(1 - F(R_i)) + F(L_i)(1 - F(L_i))]$. Define $W = \text{diag}(w_1, \dots, w_n)$.

$$H = -X^T W X \quad (17)$$

Clearly, the hessian is negative definite. Thus, the loglikelihood is strictly convex, and a unique global maximum exists.

0.1.4 IRLS

We use Newton's algorithm to find MLE for the AFT model, using negative loglikelihood (NLL). The Newton update is as follows:

$$\begin{aligned} \beta &= \tilde{\beta} - H^{-1} \tilde{g} \\ &= \tilde{\beta} - (X^T \tilde{W} X)^{-1} X^T \tilde{\mu} \\ &= (X^T \tilde{W} X)^{-1} ((X^T \tilde{W} X) \tilde{\beta} - X^T \tilde{\mu}) \\ &= (X^T \tilde{W} X)^{-1} X^T (\tilde{W} X \tilde{\beta} - \tilde{\mu}) \\ &= (X^T \tilde{W} X)^{-1} X^T (\tilde{W} X \tilde{\beta} - \tilde{\mu}) \end{aligned} \quad (18)$$

where, define the working response $\tilde{z} = X\tilde{\beta} - \widetilde{W}^{-1}\tilde{\mu}$. Here, the tilde denotes that the respective values are evaluated using the parameters from the previous step.

Hence, at each step we are solving a weighted least squares problem, which is a minimizer of:

$$\sum_{i=1}^n \tilde{w}_i (\tilde{z}_i - \tilde{x}_i^T \beta)^2 \quad (19)$$

We can rewrite z as:

$$\begin{aligned} \tilde{z}_i &= x_i^T \tilde{\beta} + \frac{\tilde{\mu}_i}{\tilde{w}_i} \\ &= x_i^T \tilde{\beta} + \end{aligned} \quad (20)$$

This algorithm is the iteratively reweighted least squares (IRLS), since at each iteration we solve a weighted least squares problem.

0.1.5 Elastic net penalty and coordinate descent

We define the elastic net (L1 + L2) penalty as follows:

$$\lambda P_\alpha(\beta) = \lambda(\alpha \|\beta\|_1 + 1/2(1 - \alpha) \|\beta\|_2^2) \quad (21)$$

Adding the elastic net (L1 + L2) penalty, we get the following penalized weighted least squares objective:

$$M = \sum_{i=1}^n \tilde{w}_i (\tilde{z}_i - \tilde{x}_i^T \beta)^2 + \lambda P_\alpha(\beta) \quad (22)$$

The subderivative of the optimization objective is given as:

$$\frac{\partial f}{\partial w_j} = \begin{cases} \{-\frac{\partial l(\beta)}{\partial w_j} + \lambda_2 w_j - \lambda_1\} & \text{if } w_j < 0 \\ [-\frac{\partial l(\beta)}{\partial w_j} - \lambda_1, -\frac{\partial l(\beta)}{\partial w_j} + \lambda_1] & \text{if } w_j = 0 \\ \{-\frac{\partial l(\beta)}{\partial w_j} + \lambda_2 w_j + \lambda_1\} & \text{if } w_j > 0 \end{cases} \quad (23)$$

Using the subderivative, upto three cases of solutions for w_j may be obtained. The coordinate descent algorithm works by cycling through each w_j in turn, keeping the others constant, and using the above estimate to calculate the optimal value β_j^* . This is repeated until convergence.

Sources:

- Machine Learning: A Probabilistic Perspective by Kevin Murphy
- glmnet and coxnet papers
- AFT: TD Hocking