# Regression in Survival analysis: The Cox model

## Motivation

Aim: study, say, survival conditional on covariates $X$,

$$S(t \mid X = x).$$

## Cox model

The Cox proportional hazards model parameterizes the hazard function as

$$\lambda(t \mid \mathbf{x}_i) = \lambda_0(t) \, e^{\beta^T \mathbf{x}_i(t)},$$

where $\beta$ is a vector of coefficients and $\mathbf{x}_i(t) = (x_{i1}(t), x_{i2}(t), \ldots, x_{ip}(t))$ are fixed or time-varying, predictable covariates.

Remember:

A statistical model $P = \{ P_\eta : \eta \in \Gamma \}$.

Parametric model $P = \{ P_\theta : \theta \in \Theta \}$ $\Theta \subseteq \mathbb{R}^k$ for a positive finite $k$.

---

Cox-model is semi-parametric because

- $\alpha_0(t)$ is non parametric, and
- $e^{\beta^T x_i(t)}$ is parametric.

---

Implications of the Cox model:

Def: the hazard ratio of two hazards $\alpha_1(t)$ and $\alpha_2(t)$ is $\dfrac{\alpha_1(t)}{\alpha_2(t)}$ $\left( \neq \dfrac{1-S_1(t)}{1-S_2(t)} \neq \dfrac{S_1(t)}{S_2(t)} \right)$

"HR"

The HR in a Cox model is $\dfrac{\alpha(t|x_1)}{\alpha(t|x_2)} = e^{\beta^T(X_1(t) - X_2(t))}$

often assumed that $X(t) = X$.

Suppose two vectors $X_1$ and $X_2$ satisfy $X_{1,j} - X_{2,j} = 1$, $X_{1,i} = X_{2,i}$ $i \neq j$, then the HR is $e^{\beta_j}$.

# Lecture 9

- Cox regression model
- Interpretation of regression coefficients, hazard ratios
- Collapsibility.

## Partial likelihood

AIM (intuitively) get rid of $\alpha_0(t)$ and only consider $\beta$.

Construction:

$$\lambda_i(t) = Z_i(t)\, \alpha_0(t)\, e^{\beta^T x_i(t)} = \alpha_0(t)\left[ Z_i(t)\, e^{\beta^T x_i(t)} \right] \qquad x_i(t) \text{ is predictable.}$$

$$N_\bullet(t) = \sum_{\ell=1}^{n} N_\ell(t), \quad \lambda_\bullet(t) = \sum_{\ell=1}^{n} \lambda_\ell(t) = \sum_{\ell=1}^{n} Z_\ell(t)\, \alpha_0(t)\, e^{\beta^T x_\ell(t)}$$

Now, $\quad \lambda_i(t) = \lambda_\bullet(t)\, \pi(i|t), \quad$ where

$$\pi(i|t) = \frac{\lambda_i(t)}{\lambda_\bullet(t)} = \frac{Z_i(t)\, \alpha_0(t)\, e^{\beta^T x_i(t)}}{\sum_{\ell=1}^{n} Z_\ell(t)\, \alpha_0(t)\, e^{\beta^T x_\ell(t)}}$$

## Def. Partial likelihood

Consider event times $T_1 < T_2 < T_3 \dots$ and let $i_j$ be the index of the individual who has an event at $T_j$. The partial likelihood is:

$$\mathcal{L}(\beta) = \prod_{T_j} \pi(i_j | T_j) = \prod_{T_j} \left( \frac{Z_{i_j}(T_j) \, e^{\beta^T x_{i_j}(T_j)}}{\sum_{l=1}^{n} Z_l(T_j) \, e^{\beta^T x_l(T_j)}} \right) = \prod_{T_j} \frac{e^{\beta^T x_{i_j}(T_j)}}{\sum_{l \in R_j} e^{\beta^T x_l(T_j)}}$$

$$R_j = \{ l : Z_i(T_j) = 1 \}$$

## Result:

Let $\hat{\beta}$ be the value of $\beta$ that maximizes $\mathcal{L}(\beta)$, called the "maximum partial likelihood estimator".

In large samples: $\hat{\beta} \sim N\left( \beta_0, \, \mathcal{I}(\beta_0)^{-1} \right)$, $\mathcal{I} = -\frac{\delta^2}{\delta \beta_h, \delta \beta_l} \log \mathcal{L}(\beta)$

"partial information matrix".

Proof technique:
Show that $U(\beta) = \frac{\delta \log \mathcal{L}(\beta)}{\delta \beta}$
is a martingale.

## From a Cox model to survival:

$$\lambda_\bullet(t) = \sum_{\ell=1}^{r} \lambda_\ell(t) = \sum_{\ell=1}^{n} Z_\ell(t)\, d_0(t)\, e^{\beta^T x(t)}$$

Survival $N_\bullet(t)$ has intensity $\lambda_\bullet(t)$.

If $\beta$ were known, we could consider a N-A "type" of estimator:

$$\hat{H}(t;\beta) = \int_0^t \frac{dN_\bullet(u)}{\sum_{\ell=1}^{n} Z_i(u)\, e^{\beta^T x_i(u)}}$$

## Breslow estimator:

$$\hat{H}_0(t;\hat{\beta}) = \int_0^t \frac{dN_\bullet(u)}{\sum_{\ell=1}^{n} Z_i(u)\, e^{\hat{\beta}^T x_i(u)}} = \sum_{T_j \le t} \frac{1}{\sum_{\ell \in R_j} e^{\hat{\beta}^T x_i(T_j)}}$$

Thus,

$$\hat{H}(t;x_0) = \hat{H}_0(t)\, e^{\hat{\beta}^T x_0},$$

$$S(t|x_0) = \prod_{u \le t}\left(1 - dH(u|x_0)\right),$$

$$\hat{S}(t|x_0) = \prod_{T_j \le t}\left(1 - \Delta\hat{H}(T_j|x_0)\right)$$

# Model checking

Consider a cox model with fixed covariates.

$$\alpha(t|x) = \alpha_0(t) e^{\beta^T x}$$

$$-\log(S(t|x)) = \int_0^t \alpha(s|x)\,ds = \int_0^t \alpha_0(s) e^{\beta^T x}\,ds$$

$$\log\left(-\log\left(S(t|x)\right)\right) = \log\left\{\int_0^t \alpha_0(s)\,ds\right\} + \beta^T x.$$

$$\log\left(-\log\left(S(t|x_1)\right)\right) - \log\left(-\log\left(S(t|x_2)\right)\right) = \beta^T x_1 - \beta^T x_2.$$

Ex: Collapsibility

Suppose I have estimates of $P(T^{a=1} > t \mid V=v)$ and $P(T^{a=0} > t \mid V=v)$ for a covariate vector $V$.

$$P(T^{a=1} > t) - P(T^{a=0} > t) = \sum_U \left[ P(T^{a=1} > t \mid U=v) - P(T^{a=0} > t \mid U=v) \right] P(V=v)$$

Survival difference is collapsible.

$$\sum_U P(V=v) = 1.$$

# Hazard ratio and collapsibility

$A \in \{0,1\}$, $Z \in [0,\infty)$ baseline covariate.

$r > 0$, $\alpha(t) > 0$ $\forall t > 0$.

Laplace transform: $L(c) = \mathbb{E}\left(e^{-cZ}\right)$ for $c \in \mathbb{C}$.

$\alpha(t \mid A=0, Z) = Z\alpha(t)$

$\alpha(t \mid A=1, Z) = rZ\alpha(t)$

$S(t \mid A=0) = L\left(H(t)\right)$, $\quad H(t) = \int_0^t \alpha(s)\,ds$.

$\alpha(t \mid A=0) = -\alpha(t)\dfrac{L'(H(t))}{L(H(t))}$.

$S(t \mid A=1) = L\left(rH(t)\right)$

$Z$ is gamma distributed with mean 1 and variance $\delta$. Then,

$$L(c) = \{1 + \delta c\}^{\frac{-1}{\delta}},$$

$$S(t) := \{1 + \delta H(t)\}^{\frac{-1}{\delta}}$$

$$\alpha(t \mid A=0) = \frac{\alpha(t)}{1 + \delta H(t)}$$

$$\alpha(t \mid A=1) = \frac{r\alpha(t)}{1 + r\delta H(t)}$$

Suppose $\delta = 1$.

$$\frac{\alpha(t \mid A=1)}{\alpha(t \mid A=0)} = r\,\frac{1 + H(t)}{1 + rH(t)} \neq r$$

for $r \neq 1$ at all $t > 0$.