

Exact gradient and Hessian for TOPALS model fitting with age-grouped input data

Carl Schmertmann

October 2021

1 Main Idea

Here I show the derivation for the analytical gradient and Hessian of the penalized likelihood with respect to TOPALS model parameters. These expressions allow optimization by Newton-Raphson iteration, which is slightly different from the iteratively reweighted least squares approach. In particular, the covariance matrix estimated from the IRLS approach at the final estimate isn't *quite* the same as the approximation that uses the inverse of the negative of the Hessian.

1.1 Objective

We wish to fit a TOPALS model mortality schedule for A single years of age $x = 0, 1, \dots, (A - 1)$ from exposure and death data for $G \leq A$ age groups: $N_1 \dots N_G$ and deaths $D_1 \dots D_G$. We include a very small roughness penalty, which will affect the fit only when there are age groups with zero exposure.

1.2 Sample Data

Observed data consists of deaths D_g and exposure N_g for closed, non-overlapping age groups $g = 1 \dots G$. Age groups are defined by a vector of $G + 1$ boundary ages – e.g. if the boundary ages are $0, 1, 5, 10, \dots, 85, 90$ then the age groups are $[0, 1), [1, 5), \dots, [85, 90)$. Denote X_g as the set of integer ages that belong to group g , and n_g as the number of ages in X_g .

Typical age groups have boundaries like $0, 1, 5, 10, \dots, 85, 90$. However, if single-year death and exposure data is available then the "groups" could be simply integer ages, in which case boundaries for A single-year "groups" would be $0, 1, 2, \dots, A$.

1.3 Poisson Log Likelihood for Grouped Data

Regardless of the age grouping in the data, we assume that there is a latent schedule for A single-year ages. Deaths at each single-year age have a Poisson

distribution with expected value equal to the (possibly unobserved) exposure N_x times the mortality rate:

$$D_x \sim \text{Poisson}(N_x \mu_x) \quad x = 0, 1, 2, \dots (A-1)$$

Under the standard assumption that age-specific deaths are statistically independent, this implies that deaths are also Poisson distributed within each age group:

$$D_g \sim \text{Poisson}\left(\sum_{x \in X_g} N_x \mu_x\right) \quad g = 1, 2, \dots G$$

Expressing this in terms of the total *observed* exposure in each age group,

$$D_g \sim \text{Poisson}\left(N_g \sum_{x \in X_g} \left[\frac{N_x}{N_g}\right] \mu_x\right)$$

or

$$D_g \sim \text{Poisson}(N_g M_g)$$

where M_g represents the exposure-weighted average mortality rate across ages in group g .

In the absence of single-year exposure data, we approximate with $M_g = \left[\frac{1}{n_g}\right] \sum_{x \in X_g} \mu_x$.¹

The entire vector of G averaged mortality rates by age group is therefore

$$\mathbf{M} = \begin{pmatrix} n_1^{-1} \dots n_1^{-1} & \dots & 0 \dots 0 \\ \vdots & \ddots & \vdots \\ 0 \dots 0 & \dots & n_G^{-1} \dots n_G^{-1} \end{pmatrix} \begin{pmatrix} \mu_0 \\ \vdots \\ \mu_{A-1} \end{pmatrix} = \mathbf{W} \boldsymbol{\mu}$$

where \mathbf{W} is a $G \times A$ matrix of weights with row sums all equal to one.²

The sample log likelihood for a single-year mortality schedule $\boldsymbol{\mu}$ is thus

$$\begin{aligned} \ln L(\boldsymbol{\mu}) &= c + \sum_{g=1}^G (D_g \ln M_g - N_g M_g) \\ &= c + \sum_{g=1}^G (D_g \ln M_g - \hat{D}_g) \end{aligned}$$

where $\hat{D}_g = N_g M_g$ represents the expected number of deaths in group g . It will also be useful to express the likelihood in matrix terms:

$$\ln L = c + [\ln M_1 \dots \ln M_G] \mathbf{D} - \hat{\mathbf{D}}' \mathbf{1} \quad (1)$$

where \mathbf{D} and $\hat{\mathbf{D}}$ are $G \times 1$ vectors of observed and expected deaths, ordered by age group, and $\mathbf{1}$ is a $G \times 1$ vector of ones.

¹An alternative approach might interpolate single-year exposure $(N_0, N_1, \dots N_{A-1})$ from the available N_g to construct more nuanced weights for intra-group averaging of mortality rates. In most populations this would have only minor effects on estimates.

²When "groups" correspond to single years, $\mathbf{W} = \mathbf{I}_A$ and $\mathbf{M} = \boldsymbol{\mu}$.

1.4 TOPALS Model Mortality Schedule

The TOPALS model for single-year mortality rates is

$$\boldsymbol{\lambda} = \boldsymbol{\lambda}^* + \mathbf{B}\boldsymbol{\alpha}$$

where $\boldsymbol{\lambda}$ is an $A \times 1$ vector of age-specific log mortality rates for ages $0 \dots (A - 1)$, $\boldsymbol{\lambda}^*$ is a vector with fixed constants representing a standard schedule for those ages, \mathbf{B} is a $A \times K$ matrix of linear B-spline constants³, and $\boldsymbol{\alpha}$ is a K -dimensional vector representing deviations from the standard log mortality schedule at specified ages.

In this model the log mortality rate at age $x \in \{0 \dots (A - 1)\}$ is

$$\ln \mu_x = \lambda_x = \lambda_x^* + \mathbf{b}'_x \boldsymbol{\alpha}$$

and the mortality rate is

$$\mu_x = \exp(\lambda_x^* + \mathbf{b}'_x \boldsymbol{\alpha})$$

where \mathbf{b}'_x is the $1 \times K$ row of \mathbf{B} that corresponds to age x .

1.4.1 TOPALS derivatives

Derivatives of age-specific rates with respect to TOPALS parameters $\boldsymbol{\alpha}$ are $K \times 1$ vectors:

$$\frac{\partial \ln \mu_x}{\partial \boldsymbol{\alpha}} = \mathbf{b}_x \quad , \quad \frac{\partial \mu_x}{\partial \boldsymbol{\alpha}} = \mu_x \mathbf{b}_x$$

For the entire $G \times 1$ vector of group mortality rates, this implies

$$\begin{aligned} \frac{\partial \mathbf{M}'}{\partial \boldsymbol{\alpha}} &= \frac{\partial}{\partial \boldsymbol{\alpha}} (\boldsymbol{\mu}' \mathbf{W}') \\ &= \left(\frac{\partial \mu_0}{\partial \boldsymbol{\alpha}} \quad \dots \quad \frac{\partial \mu_{A-1}}{\partial \boldsymbol{\alpha}} \right) \mathbf{W}' \\ &= (\mu_0 \mathbf{b}_0 \quad \dots \quad \mu_{A-1} \mathbf{b}_{A-1}) \mathbf{W}' \\ &= \mathbf{B}' \text{diag}(\boldsymbol{\mu}) \mathbf{W}' \end{aligned}$$

Abbreviate this as

$$\frac{\partial \mathbf{M}'}{\partial \boldsymbol{\alpha}} = \mathbf{X}'$$

remembering that the $G \times K$ matrix $\mathbf{X} = \mathbf{W} \text{diag}(\boldsymbol{\mu}) \mathbf{B}$ varies with parameters $\boldsymbol{\alpha}$ via the $\boldsymbol{\mu}$ terms in the central diagonal matrix.

³Typically $K = 7$, with spline knots fixed at ages $t = 0, 1, 10, 20, 40, 70, 99$. Values in the k th column of \mathbf{B} are

$$B_{xk} = \begin{cases} \frac{x - t_{k-1}}{t_k - t_{k-1}} & \text{for } x \in [t_{k-1}, t_k] \\ \frac{t_{k+1} - x}{t_{k+1} - t_k} & \text{for } x \in [t_k, t_{k+1}] \\ 0 & \text{otherwise} \end{cases}$$

The $G \times 1$ vector of expected deaths in groups $1 \dots G$ is

$$\hat{\mathbf{D}} = \begin{pmatrix} N_1 M_1 \\ \vdots \\ N_G M_G \end{pmatrix} = \text{diag}(\mathbf{N}) \mathbf{M}$$

and the derivative of its transpose with respect to TOPALS parameters $\boldsymbol{\alpha}$ is a $K \times G$ matrix:

$$\frac{\partial \hat{\mathbf{D}}'}{\partial \boldsymbol{\alpha}} = \frac{\partial \mathbf{M}'}{\partial \boldsymbol{\alpha}} \text{diag}(\mathbf{N}) = \mathbf{X}' \text{diag}(\mathbf{N}) \quad (2)$$

1.4.2 Penalized Log Likelihood for TOPALS parameters

With the TOPALS parameterization the log likelihood for a sample $\{D_g, N_g\}$ is

$$\begin{aligned} \ln L(\boldsymbol{\alpha}) &= c + \sum_{g=1}^G \left(D_g \ln M_g(\boldsymbol{\alpha}) - \hat{D}_g(\boldsymbol{\alpha}) \right) \\ &= c + [\ln M_1 \dots \ln M_G] \mathbf{D} - \hat{\mathbf{D}}' \mathbf{1} \end{aligned}$$

To stabilize estimates in small populations with few deaths, we add a small penalty term to the log likelihood:

$$\begin{aligned} Q(\boldsymbol{\alpha}) &= \ln L(\boldsymbol{\alpha}) - \text{penalty}(\boldsymbol{\alpha}) \\ &= c + [\ln M_1 \dots \ln M_G] \mathbf{D} - \hat{\mathbf{D}}' \mathbf{1} - \frac{1}{2} \boldsymbol{\alpha}' \mathbf{P} \boldsymbol{\alpha} \end{aligned} \quad (3)$$

where \mathbf{P} is a $K \times K$ matrix of constants selected so that the penalty term equals the sum of squared differences between consecutive α parameters – i.e., $\frac{1}{2} \boldsymbol{\alpha}' \mathbf{P} \boldsymbol{\alpha} = \sum_{k=2}^K (\alpha_k - \alpha_{k-1})^2$. Adding the penalty gives priority to sets of TOPALS parameters $(\alpha_1 \dots \alpha_K)$ that are similar to one another, and thus to log mortality schedules that look more like up-and-down vertical shifts of the standard schedule. For all but the smallest populations the penalty term has virtually no effect on parameter estimates. For very small populations with zero deaths in some age groups, the addition of the penalty stabilizes estimated schedules by borrowing strength across groups.

1.4.3 Gradient

We want to select $\boldsymbol{\alpha}$ to maximize Q . This requires setting a vector of derivatives equal to zero: $\frac{\partial Q}{\partial \boldsymbol{\alpha}} = 0 \in \mathbb{R}^K$.

Differentiating Eq. (3) with respect to the TOPALS parameters and substituting some of the results above produces an analytical expression for the $K \times 1$ gradient vector of first derivatives:

$$\begin{aligned}
g(\alpha) = \frac{\partial Q}{\partial \alpha} &= \left[\frac{1}{M_1} \frac{\partial M_1}{\partial \alpha} \dots \frac{1}{M_G} \frac{\partial M_G}{\partial \alpha} \right] \mathbf{D} - \frac{\partial \hat{\mathbf{D}}'}{\partial \alpha} \mathbf{1} - \mathbf{P}\alpha \\
&= \frac{\partial \mathbf{M}'}{\partial \alpha} \text{diag}\left(\frac{1}{\mathbf{M}}\right) \mathbf{D} - \frac{\partial \mathbf{M}'}{\partial \alpha} \text{diag}(\mathbf{N}) \mathbf{1} - \mathbf{P}\alpha \\
&= \mathbf{X}' \left[\text{diag}\left(\frac{1}{\mathbf{M}}\right) \mathbf{D} - \mathbf{N} \right] - \mathbf{P}\alpha
\end{aligned} \tag{4}$$

1.4.4 Hessian

In Equation (4) the terms that vary with TOPALS coefficients α are \mathbf{X} , $\text{diag}(\frac{1}{\mathbf{M}})$, and α . In order to construct the Hessian matrix of second derivatives, start by considering how the $K \times 1$ gradient changes with a change in only one of the α s – say, α_3 .

By the chain rule,

$$\begin{aligned}
\frac{\partial}{\partial \alpha_3} \left(\frac{\partial Q}{\partial \alpha} \right) &= \left[\frac{\partial \mathbf{X}'}{\partial \alpha_3} \right] \left[\text{diag}\left(\frac{1}{\mathbf{M}}\right) \mathbf{D} - \mathbf{N} \right] \\
&\quad + \mathbf{X}' \frac{\partial}{\partial \alpha_3} \left[\text{diag}\left(\frac{1}{\mathbf{M}}\right) \right] \mathbf{D} \\
&\quad - \frac{\partial}{\partial \alpha_3} [\mathbf{P}\alpha]
\end{aligned} \tag{5}$$

Equation (5) includes several matrix derivatives:

$$\begin{aligned}
\underbrace{\left[\frac{\partial \mathbf{X}'}{\partial \alpha_3} \right]}_{K \times G} &= \mathbf{B}' \left[\text{diag} \left(\frac{\partial \mu}{\partial \alpha_3} \right) \right] \mathbf{W}' \\
&= \mathbf{B}' \left[\text{diag} \left(e'_3 \frac{\partial \mu'}{\partial \alpha} \right) \right] \mathbf{W}' \\
&= \mathbf{B}' [\text{diag} (b'_3 \text{diag}(\mu))] \mathbf{W}' \\
&= \mathbf{B}' [\text{diag} (b_3 \circ \mu)] \mathbf{W}'
\end{aligned} \tag{6}$$

where b_3 is the $A \times 1$ third column of \mathbf{B} and the \circ symbol represents element-by-element multiplication.

$$\begin{aligned}
\frac{\partial}{\partial \alpha_3} \left[\text{diag} \left(\frac{1}{\mathbf{M}} \right) \right] &= \text{diag} \left(\frac{\partial}{\partial \alpha_3} \left[\frac{1}{M_1} \cdots \frac{1}{M_G} \right] \right) \\
&= -\text{diag} \left(\left[\frac{\partial M_1}{\partial \alpha_3} M_1^{-2} \cdots \frac{\partial M_G}{\partial \alpha_3} M_G^{-2} \right] \right) \\
&= -\text{diag} \left(\left[\frac{\partial M_1}{\partial \alpha_3} \cdots \frac{\partial M_G}{\partial \alpha_3} \right] \text{diag} (M_1^{-2} \cdots M_G^{-2}) \right) \quad (7) \\
&= -\text{diag} ([b'_3 \text{diag}(\mu) \mathbf{W}'] \text{diag} (M_1^{-2} \cdots M_G^{-2})) \\
\text{or} \\
&= -\text{diag} \left(\frac{(b_3 \circ \mu)' w_1}{M_1^2} \cdots \frac{(b_3 \circ \mu)' w_G}{M_G^2} \right)
\end{aligned}$$

where w_g is the $A \times 1$ g -th column of \mathbf{W} .

Thus the *third* column of the Hessian matrix $Q_{\theta\theta'}$ is

$$\begin{aligned}
\frac{\partial}{\partial \alpha_3} \left(\frac{\partial Q}{\partial \alpha} \right) &= \mathbf{B}' [\text{diag} (b_3 \circ \mu)] \mathbf{W}' \left[\text{diag} \left(\frac{1}{\mathbf{M}} \right) \mathbf{D} - \mathbf{N} \right] \\
&\quad - \mathbf{X}' \left[\text{diag} \left(\frac{(b_3 \circ \mu)' w_1}{M_1^2} \cdots \frac{(b_3 \circ \mu)' w_G}{M_G^2} \right) \right] \mathbf{D} \quad (8) \\
&\quad - [\mathbf{P} e_3]
\end{aligned}$$

and the generic j -th column of the Hessian is the $K \times 1$ vector

$$\begin{aligned}
h_j = \frac{\partial}{\partial \alpha_j} \left(\frac{\partial Q}{\partial \alpha} \right) &= \mathbf{B}' [\text{diag} (b_j \circ \mu)] \mathbf{W}' \left[\text{diag} \left(\frac{1}{\mathbf{M}} \right) \mathbf{D} - \mathbf{N} \right] \\
&\quad - \mathbf{X}' \left[\text{diag} \left(\frac{(b_j \circ \mu)' w_1}{M_1^2} \cdots \frac{(b_j \circ \mu)' w_G}{M_G^2} \right) \right] \mathbf{D} \quad (9) \\
&\quad - [\mathbf{P} e_j]
\end{aligned}$$

And the complete $K \times K$ Hessian matrix is the (admittedly complicated, but computable)

$$\mathbf{H}(\alpha) = \mathbf{Q}_{\theta\theta'} = [h_1 \cdots h_K] \quad (10)$$

1.5 Newton-Raphson Iteration

With analytical expressions for the gradient and Hessian [Equations (4) and (10), respectively] we can now use Newton-Raphson iteration to solve for the $\alpha \in \mathbb{R}^K$ that maximizes the penalized likelihood in Equation (3).

Starting at an arbitrary vector α_0 (usually all zeroes) for the TOPALS offsets, we start at $t = 0$ and repeat the following until convergence:

1. calculate the gradient and Hessian at the current parameter values: $g(\alpha_t)$ and $H(\alpha_t)$
2. using the (multivariate) quadratic approximation to the objective function, generate a new parameter vector at which the gradient should approximately equal zero:

$$\alpha_{t+1} = \alpha_t - [H(\alpha_t)]^{-1}g(\alpha_t)$$

3. increment t by one and return to step 1