the coefficients of $\mathbf{x}^{(j)}$ and $\mathbf{x}^{(k)}$ are not well determined individually, but their model part $(\beta_j \mathbf{x}^{(j)} + \beta_k \mathbf{x}^{(k)})$ still is. Geometrically, in the $\beta_j$-$\beta_k$-plane the high confidence region forms narrow ellipses, i.e., the $\beta$ components themselves are linearly related, or that the coefficients of $\hat{\beta}_j$ and $\hat{\beta}_k$ themselves highly correlated but not be well determined individually, i.e., have a large variance. In the extreme case of "perfect" correlation, the matrix $\mathbf{X}$ would have columns $j$ and $k$ collinear and hence only have rank $\leq p - 1$. When the correlation is less extreme, $\mathbf{X}$ is still of full rank $p$ and $\mathbf{X}^{\mathsf{T}}\mathbf{X}$ is close to a singular matrix.[2] One approach to this problem is to *regularize* it by improving the condition of the matrix corresponding to $\mathbf{X}^{\mathsf{T}}\mathbf{X}$.

To give a numerical example, say $\mathbf{x}^{(2)} \approx -2\mathbf{x}^{(1)}$, then $1\mathbf{x}^{(1)}$ is close to $3\mathbf{x}^{(1)} + \mathbf{x}^{(2)}$ and hence to $5\mathbf{x}^{(1)} + 2\mathbf{x}^{(2)}$ or $51\mathbf{x}^{(1)} + 25\mathbf{x}^{(2)}$. One way to make the linear combination more clearly determined is to restrict the coefficients to small (absolute) values, or, more conveniently requiring that $\sum_j \beta_j^2$ be "small". This leads to the so called *ridge regression*

$$\tilde{\boldsymbol{\beta}}(s) = \underset{\|\boldsymbol{\beta}\|^2 \leq s}{\arg\min} \|\mathbf{Y} - X\boldsymbol{\beta}\|^2,$$

which can be shown (by way of a Lagrange multiplier) to be equivalent to

$$\widehat{\boldsymbol{\beta}}^*(\lambda) = \underset{\boldsymbol{\beta}}{\arg\min} \{\|\mathbf{Y} - X\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|^2\}, \tag{7.7}$$

where there is a one-to-one relationship between $\lambda$ and the bound $s$ above, and the limit $\lambda \to 0$ corresponds to $s = \max \|\boldsymbol{\beta}\|^2 \to \infty$, namely the ordinary least squares of chapter 1. As there, by setting derivatives $\partial/\partial\beta$ to zero, this minimization problem is equivalent to the "normal equations"

$$(X^{\mathsf{T}}X + \lambda I)\widehat{\boldsymbol{\beta}}^* = X^{\mathsf{T}}\mathbf{Y}, \tag{7.8}$$

where the $p \times p$ matrix $(X^{\mathsf{T}}X + \lambda I)$ will be non-singular (and "well-conditioned") as soon as $\lambda > 0$ is large enough, even when $n < p$ and $X^{\mathsf{T}}X$ is clearly singular.

The ridge penalty entails that $\hat{\beta}_j(\lambda) \to 0$ ( *"shrinking"*) when $\lambda \to \infty$, and also, in general, $\hat{\beta}_j \to \hat{\beta}_{j'}$ ("shrinking together") for two different coefficients.

Therefore, it's intuitive that $\widehat{\boldsymbol{\beta}}^*$ will have some bias ($\mathbb{E}[\widehat{\boldsymbol{\beta}}^*] \neq \boldsymbol{\beta}$), but that its variance(s) can be considerably smaller than $\widehat{\boldsymbol{\beta}}^{LS}$ such that mean squared errors are smaller. As for smoothing (in particular, spline smoothing, sect. 3.4) we have a regularization parameter $\lambda$ which determines the trade-off between bias and variance, and as there, we'd use something like cross validation to determine an approximately optimal value for $\lambda$.

In the literature (and the R function `lm.ridge()` from the package MASS) there are "cheaper" approaches like GCV (see ch. 4) for determining an approximately optimal $\lambda$. In practice, one often wants to look at the *ridge traces*, i.e., a plot of the coefficients $\hat{\beta}_j(\lambda)$ vs $\lambda$. As an example we consider the longley macroeconomical data, for once modelling $y = $ `GNP.deflator` as function of the other six variables. The ridge traces $\hat{\beta}_j(\lambda)$ are shown in Figure 7.11. We have used a "relevant" interval for $\lambda$ where the shrinking towards zero is visible (but still somewhat distant from the limit).

### 7.7.3 The Lasso

In some sense, the "Lasso" (Tibshirani, 1996) regression is just a simple variant of ridge regression. However with the goal of *variable selection* in mind, the lasso brings a major improvement:

---

[2]such that $\mathrm{Cov}(\widehat{\boldsymbol{\beta}}) = \sigma^2(X^{\mathsf{T}}X)^{-1}$ (section 1.4.1) will have very large entries corresponding to the high variance (and correlation) of $\hat{\beta}_j$ and $\hat{\beta}_k$.
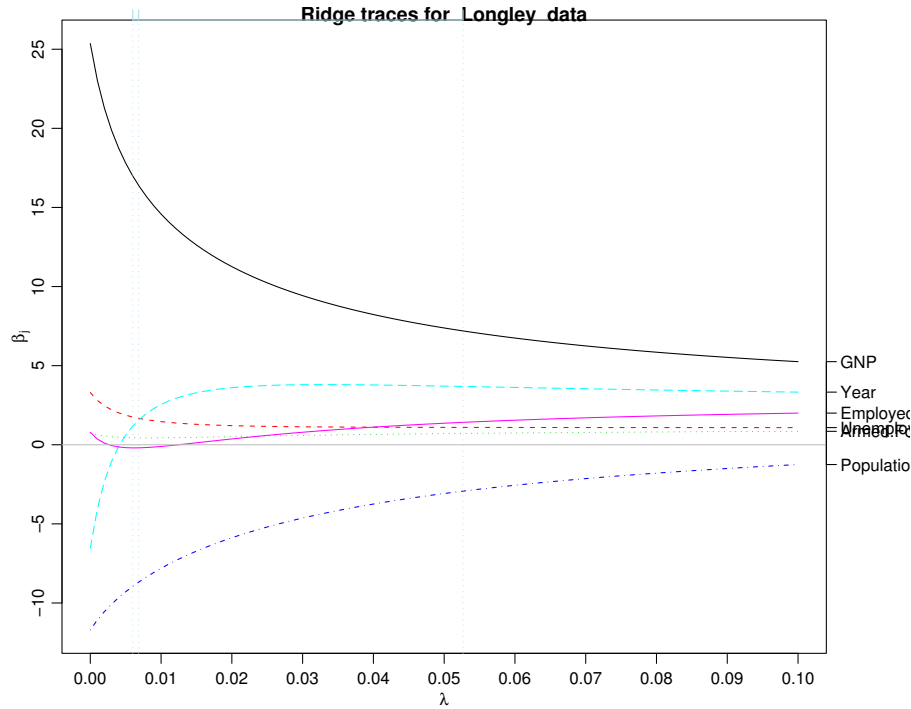
Figure 7.11: Ridge traces for the six coefficients $\beta_j(\lambda)$ $(j = 1, \ldots, 6)$ for the Longley data. The vertical lines indicate traditional estimates of the optimal ridge parameter $\lambda$.

The lasso can be defined by restricting the *absolute* instead of the squared values of the coefficients, i.e.,

$$\tilde{\boldsymbol{\beta}}(s) = \underset{\sum_j |\beta_j| \leq s}{\arg\min} \|\mathbf{Y} - X\boldsymbol{\beta}\|^2,$$

or

$$
\begin{aligned}
\widehat{\boldsymbol{\beta}}^*(\lambda) &= \underset{\boldsymbol{\beta}}{\arg\min} \ \{\|\mathbf{Y} - X\boldsymbol{\beta}\|^2 + \lambda \sum_{j=1}^{p} |\beta_j|\}, \\
&= \underset{\boldsymbol{\beta}}{\arg\min} \ \{\|\mathbf{Y} - X\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|_1\}.
\end{aligned}
\tag{7.9}
$$

As opposed to the ridge regression case above, this problem is not solvable by simple linear algebra but rather needs quadratic programming or related algorithms.

On the other hand, the solution is much more interesting, because it will be frequent that $\hat{\beta}_j$ will become *exactly 0* as soon as $\lambda > \lambda_j$, in other words, choosing $\lambda$ here, automatically means model selection, namely only choosing regressor variables $\mathbf{x}^{(j)}$ with $\beta_j \neq 0$.

This can be visualized considering the "lasso traces".