Eigenvalue based regression, including Ridge (Tikhonov) and Differential regularization

Antoine Savine, April 2020

The exercise is to perform a classic least square linear regression $Y = \mu_Y + (\phi - \mu_\phi)\beta$, where the columns of $\phi = \phi(X)$ are basis functions (e.g. monomials, excluding constant) of known inputs $X$ (also excluding constant, with examples in rows and inputs in columns), given a column vector $Y$ of the corresponding targets, where $\mu_Y$ is the mean of $Y$ and the row vector $\mu_\phi$ contains the means of the columns of $\phi$. To simplify notations, we denote $\phi \equiv \phi - \mu_\phi$ and $Y \equiv Y - \mu_Y$. Classic least squares finds $\beta$ by minimization of the least square errors:

$$\beta = \arg\min_\beta \|Y - \phi\beta\|^2$$

The analytic solution, also called 'normal equation':

$$\beta = (\phi^T \phi)^{-1} \phi^T Y$$

is known to bear unstable results, the matrix $\phi^T \phi$ often being near singular (certainly so with monomials of high degree of correlated inputs). This is usually resolved with SVD regression. We prefer the (very similar) eigenvalue regression, which we recall first, and then, extend to Ridge (Tikhonov) regularization and differential regression (in the sense of Huge and Savine, Risk 2020).

Parts 1 and 2 are summaries of classic results. Part 3 is new.

After $\beta$ is learned, the value approximation for an input row vector $x$ is given by $y = \mu_Y + [\phi(x) - \mu_\phi]\beta$ and the derivative approximations are given by:

$$y_j = \phi_j(x)\beta$$

where subscripts denote partial derivatives to input number j.


1. Standard eigenvalue regression

   Perform the eigenvalue decomposition $\phi^T \phi = PDP^T$.


   Denote $D^{-\frac{1}{2}}$ the diagonal matrix whole diagonal elements are the elements of the diagonal matrix $D$, raised to power $-0.5$, when they exceed a threshold (say, $10^{-8}$ times the mean trace of $D$), and zero otherwise.

Denote $\phi = \phi PD^{-\frac{1}{2}}$ and perform the least square minimization in the orthonormal basis:

$$\beta = \arg\min_\beta \left\| Y - \phi\beta \right\|^2$$

The normal equation is stable in the orthonormal basis:

$$\beta = \left( \phi^T \phi \right)^{-1} \phi^T Y$$

It is easy to see that $\phi^T \phi$ is a diagonal matrix with diagonal elements 1 corresponding to significant eigenvalues, and 0 corresponding to insignificant ones. With the convention $\left( \phi^T \phi \right)^{-1} = \phi^T \phi$ (invert the significant diagonal elements and zero the insignificant ones), we get:

$$\beta = D^{-\frac{1}{2}} P^T \phi^T Y$$

(notice, $D^{-\frac{1}{2}}$ zeroes the lines corresponding to insignificant eigenvalues so there is no need to left multiply by $\phi^T \phi$.)

Hence: $Y = \phi\beta = \phi PD^{-1}P^T \phi^T Y = \beta Y$ where $D^{-1} = \left( D^{-\frac{1}{2}} \right)^2$ has diagonal elements inverse of the significant eigenvalues in $D$, zero for insignificant eigenvalue, and:

$$\boxed{\beta = PD^{-1}P^T \phi^T Y}$$

2. Ridge (Tikhonov) eigenvalue regression

Classic regression works best with regularization, the most common classic form of which is ridge regression, also called Tikhonov regularization, which adds a penalty on the norm of $\beta$ in the objective cost.

$$\beta = \arg\min_\beta \left[ \left\| Y - \phi\beta \right\|^2 + \lambda^2 \left\| \beta \right\|^2 \right] = \arg\min_\beta \left[ \left\| Y - \left( \phi PD^{-\frac{1}{2}} \right)\left( D^{\frac{1}{2}}P^T \beta \right) \right\|^2 + \lambda^2 \left\| \beta \right\|^2 \right]$$

$$= PD^{-\frac{1}{2}} \arg\min_\gamma \left[ \left\| Y - \phi\gamma \right\|^2 + \lambda^2 \left\| PD^{-\frac{1}{2}}\gamma \right\|^2 \right] = PD^{-\frac{1}{2}} \arg\min_\gamma \left[ \left\| Y - \phi\gamma \right\|^2 + \lambda^2 \left\| D^{-\frac{1}{2}}\gamma \right\|^2 \right]$$

$$= PD^{-\frac{1}{2}} \left[ \phi^T \phi + \lambda^2 D^{-1} \right]^{-1} \phi^T Y = PD^{-\frac{1}{2}} \left[ \phi^T \phi + \lambda^2 D^{-1} \right]^{-1} D^{-\frac{1}{2}}P^T \phi^T Y = P\Lambda^{-1}P^T \phi^T Y$$

where $\Lambda^{-1}$ has diagonal elements $\dfrac{1}{D_{ii} + \lambda^2}$ where $D_{ii}$ is significant, zero otherwise. And we get:

$$\beta(\lambda) = P\Lambda(\lambda)^{-1} P^T \phi^T Y$$

The Tikhonov parameter $\lambda$ can be found e.g. by cross validation:

$$\lambda = \arg\min_\lambda \left\| Y_V - \phi_V \beta(\lambda) \right\|^2$$

where $\phi_V = \phi(X_V)$, $(X_V, Y_V)$ form a validation set of independent examples and $\beta(\lambda)$ is the result of a ridge regression over the training set with Tikhonov parameter $\lambda$, obtained with the boxed formula above. The objective function can be expanded:

$$
\begin{aligned}
f(\lambda) &= \left\| Y_V - \phi_V \beta(\lambda) \right\|^2 \\
&= \left( Y_V - \phi_V \beta(\lambda) \right)^T \left( Y_V - \phi_V \beta(\lambda) \right) \\
&= \left( Y_V - \phi_V P\Lambda(\lambda)^{-1} P^T \phi^T Y \right)^T \left( Y_V - \phi_V P\Lambda(\lambda)^{-1} P^T \phi^T Y \right) \\
&= \left( Y_V^T - Y^T \phi P\Lambda(\lambda)^{-1} P^T \phi_V^T \right) \left( Y_V - \phi_V P\Lambda(\lambda)^{-1} P^T \phi^T Y \right) \\
&= Y_V^T Y_V - Y^T \phi P\Lambda(\lambda)^{-1} P^T \phi_V^T Y_V - Y_V^T \phi_V P\Lambda(\lambda)^{-1} P^T \phi^T Y \\
&\quad + Y^T \phi P\Lambda(\lambda)^{-1} P^T \phi_V^T \phi_V P\Lambda(\lambda)^{-1} P^T \phi^T Y
\end{aligned}
$$

Since $Y_V^T Y_V$ doesn't depend on $\lambda$, we minimize:

$$
\begin{aligned}
g(\lambda) &= Y^T \phi P\Lambda(\lambda)^{-1} P^T \phi_V^T \phi_V P\Lambda(\lambda)^{-1} P^T \phi^T Y \\
&\quad - Y^T \phi P\Lambda(\lambda)^{-1} P^T \phi_V^T Y_V \\
&\quad - Y_V^T \phi_V P\Lambda(\lambda)^{-1} P^T \phi^T Y \\
&= K^T \Lambda(\lambda)^{-1} M \Lambda(\lambda)^{-1} K - K^T (\lambda)^{-1} L - L^T \Lambda(\lambda)^{-1} K \\
&= K^T \Lambda(\lambda)^{-1} M \Lambda(\lambda)^{-1} K - 2K^T (\lambda)^{-1} L \\
&= K^T \Lambda(\lambda)^{-1} \left[ M \Lambda(\lambda)^{-1} K - 2L \right]
\end{aligned}
$$

where $K = P^T \phi^T Y \{n \times 1\}, L = P^T \phi_V^T Y_V \{n \times 1\}, M = P^T \phi_V^T \phi_V P \{n \times n\}$

Optimization may be efficiently performed by a classic one-dimensional minimization procedure.

3. Differential regression (Huge and Savine, Risk 2020)

In addition to inputs $X$ and labels $Y$, we have derivatives labels $Z$ whose columns $Z_j$ are the differentials of $Y$ to $X_j$. Denote $\phi_j$ the matrix of derivatives of the basis functions $\phi$ to $X_j$.

Linear regression makes value predictions $Y = \phi\beta$ and derivatives predictions $Z_j = \phi_j\beta$. We now minimize a cost combining value and derivatives errors:

$$\beta = \arg\min_\beta \left[ \|Y - \phi\beta\|^2 + \sum_j \lambda_j \|Z_j - \phi_j\beta\|^2 \right]$$

where $\lambda_j = \lambda^2 \dfrac{\|Y\|^2}{\|Z_j\|^2}$ (norms are computed across examples) ensures that the components of the cost are of similar magnitudes. The parameter $\lambda$ has little effect and generally left to 1.

It is not hard to see that this minimization is analytically solved with the adjusted normal equation:

$$\beta = \left( \phi^T\phi + \sum_j \lambda_j \phi_j^T \phi_j \right)^{-1} \left( \phi^T Y + \sum_j \lambda_j \phi_j^T Z_j \right)$$

This is, again, a theoretical equation, unstable in practice. As before, we change basis by eigenvalue decomposition of:

$$\phi^T\phi + \sum_j \lambda_j \phi_j^T \phi_j = PDP^T$$

-- beware, notations have changed so $P$ and $D$ denote different (respectively unitary and diagonal) matrices than before.

Changing basis as before: $\tilde\phi = \phi PD^{-\frac{1}{2}}$ (where, as previously, $D^{-\frac{1}{2}}$ has zero diagonal elements where the eigenvalues in $D$ are insignificant) we notice that:

$$\tilde\phi_j \equiv \frac{\partial\tilde\phi}{\partial X_j} = \frac{\partial\left(\phi PD^{-\frac{1}{2}}\right)}{\partial X_j} = \frac{\partial\phi}{\partial X_j} PD^{-\frac{1}{2}} = \phi_j PD^{-\frac{1}{2}}$$

Performing the minimization in the 'tilde' basis:

$$\tilde\beta = \arg\min_{\tilde\beta} \left[ \|Y - \tilde\phi\tilde\beta\|^2 + \sum_j \lambda_j \|Z_j - \tilde\phi_j\tilde\beta\|^2 \right]$$

We have the normal equation:

$$\beta = \left( \phi^T \phi + \sum_j \lambda_j \phi_j^T \phi_j \right)^{-1} \left( \phi^T Y + \sum_j \lambda_j \phi_j^T Z_j \right)$$

$$= \left[ \left( D^{-\frac{1}{2}} P^T \right) \left( \phi^T \phi + \sum_j \lambda_j \phi_j^T \phi_j \right) \left( P D^{-\frac{1}{2}} \right) \right]^{-1} \left[ \left( D^{-\frac{1}{2}} P^T \right) \left( \phi^T Y + \sum_j \lambda_j \phi_j^T Z_j \right) \right]$$

$$= D^{-\frac{1}{2}} P^T \left( \phi^T Y + \sum_j \lambda_j \phi_j^T Z_j \right)$$

Predicted values are given by:

$$Y = \phi \beta = \phi P D^{-\frac{1}{2}} D^{-\frac{1}{2}} P^T \left( \phi^T Y + \sum_j \lambda_j \phi_j^T Z_j \right) = \phi P D^{-1} P^T \left( \phi^T Y + \sum_j \lambda_j \phi_j^T Z_j \right) = \phi \beta$$

where $D^{-1} = \left( D^{-\frac{1}{2}} \right)^2$ is defined as previously (with zeroes on insignificant eigenvalues) and:

$$\boxed{\beta = P D^{-1} P^T \left( \phi^T Y + \sum_j \lambda_j \phi_j^T Z_j \right)}$$

Note that this is all consistent, in particular, derivatives predictions are given by:

$$Z_j = \phi_j \beta = \phi_j P D^{-1} P^T \left( \phi^T Y + \sum_j \lambda_j \phi_j^T Z_j \right) = \left( \phi_j P D^{-\frac{1}{2}} \right) \left[ D^{-\frac{1}{2}} P^T \left( \phi^T Y + \sum_j \lambda_j \phi_j^T Z_j \right) \right] = \phi_j \beta$$

Hence, just like Tikhonov regularization, differential regularization is analytic and extremely effective, as seen in the demonstration notebook. Unlike Tikhonov, differential regularization is *unbiased*, as demonstrated in an appendix of Huge & Savine, Risk 2020.