

Name: Jared Yu Date: May 31, 2018

1(c) Suppose for the moment that σ^2 is known. I want to get an unbiased estimate for the L_2 -risk:

$$\text{risk}(\lambda) = E\|\mathbf{f} - \hat{\mathbf{f}}_\lambda\|^2.$$

Let consider another expression for the risk first.

$$\begin{aligned} \text{risk}(\lambda) &= E\|\mathbf{f} - \hat{\mathbf{f}}_\lambda\|^2 \\ &= E[(\mathbf{f} - \hat{\mathbf{f}}_\lambda)^T (\mathbf{f} - \hat{\mathbf{f}}_\lambda)] \\ &= E[(\mathbf{f}^T - \mathbf{y}^T \mathbf{H}_\lambda^T)(\mathbf{f} - \mathbf{H}_\lambda \mathbf{y})] \\ &= E(\mathbf{f}^T \mathbf{f} - \mathbf{f}^T \mathbf{H}_\lambda \mathbf{y} - \mathbf{y}^T \mathbf{H}_\lambda^T \mathbf{f} + \mathbf{y}^T \mathbf{H}_\lambda^T \mathbf{H}_\lambda \mathbf{y}) \\ &= \mathbf{f}^T \mathbf{f} - \mathbf{f}^T \mathbf{H}_\lambda \mathbf{f} - \mathbf{f}^T \mathbf{H}_\lambda^T \mathbf{f} + \text{tr}[E(\mathbf{y} \mathbf{y}^T) \mathbf{H}_\lambda^T \mathbf{H}_\lambda] \\ &= \mathbf{f}^T \mathbf{f} - 2\mathbf{f}^T \mathbf{H}_\lambda \mathbf{f} + \text{tr}[(\mathbf{f} \mathbf{f}^T + \sigma^2 \mathbf{I}) \mathbf{H}_\lambda^T \mathbf{H}_\lambda] \\ &= \mathbf{f}^T \mathbf{f} - 2\mathbf{f}^T \mathbf{H}_\lambda \mathbf{f} + \mathbf{f}^T \mathbf{H}_\lambda^T \mathbf{H}_\lambda \mathbf{f} + \sigma^2 \text{tr}(\mathbf{H}_\lambda^T \mathbf{H}_\lambda) \\ &= \|\mathbf{f} - \mathbf{H}_\lambda \mathbf{f}\|^2 + \sigma^2 \text{tr}(\mathbf{H}_\lambda^T \mathbf{H}_\lambda). \end{aligned}$$

Then, I can start to look for an unbiased estimator for the risk using the following expectation.

$$\begin{aligned} E\|\mathbf{y} - \hat{\mathbf{f}}_\lambda\|^2 &= E\|\mathbf{y} - \mathbf{H}_\lambda \mathbf{y}\|^2 \\ &= \|(\mathbf{I} - \mathbf{H}_\lambda) \mathbf{f}\|^2 + E[\|\mathbf{y} - \mathbf{H}_\lambda \mathbf{y}\|^2 - \|(\mathbf{I} - \mathbf{H}_\lambda) \mathbf{f}\|^2] \\ &= \|(\mathbf{I} - \mathbf{H}_\lambda) \mathbf{f}\|^2 + E[\mathbf{y}^T (\mathbf{I} - \mathbf{H}_\lambda)^T (\mathbf{I} - \mathbf{H}_\lambda) \mathbf{y} - \mathbf{f}^T (\mathbf{I} - \mathbf{H}_\lambda)^T (\mathbf{I} - \mathbf{H}_\lambda) \mathbf{f}] \\ &= \|(\mathbf{I} - \mathbf{H}_\lambda) \mathbf{f}\|^2 + E\{\text{tr}[\mathbf{y}^T (\mathbf{I} - \mathbf{H}_\lambda)^T (\mathbf{I} - \mathbf{H}_\lambda) \mathbf{y}] - \text{tr}[\mathbf{f}^T (\mathbf{I} - \mathbf{H}_\lambda)^T (\mathbf{I} - \mathbf{H}_\lambda) \mathbf{f}]\} \\ &= \|(\mathbf{I} - \mathbf{H}_\lambda) \mathbf{f}\|^2 + E\{\text{tr}[\mathbf{y} \mathbf{y}^T (\mathbf{I} - \mathbf{H}_\lambda)^T (\mathbf{I} - \mathbf{H}_\lambda)] - \text{tr}[\mathbf{f} \mathbf{f}^T (\mathbf{I} - \mathbf{H}_\lambda)^T (\mathbf{I} - \mathbf{H}_\lambda)]\} \\ &= \|(\mathbf{I} - \mathbf{H}_\lambda) \mathbf{f}\|^2 + \text{tr}\{E(\mathbf{y} \mathbf{y}^T - \mathbf{f} \mathbf{f}^T) [(\mathbf{I} - \mathbf{H}_\lambda)^T (\mathbf{I} - \mathbf{H}_\lambda)]\} \\ &= \|(\mathbf{I} - \mathbf{H}_\lambda) \mathbf{f}\|^2 + \text{tr}\{\sigma^2 [(\mathbf{I} - \mathbf{H}_\lambda)^T (\mathbf{I} - \mathbf{H}_\lambda)]\} \\ &= \|(\mathbf{I} - \mathbf{H}_\lambda) \mathbf{f}\|^2 + \sigma^2 \{\text{tr}(\mathbf{I} - \mathbf{H}_\lambda^T - \mathbf{H}_\lambda + \mathbf{H}_\lambda \mathbf{H}_\lambda^T)\} \\ &= \|\mathbf{f} - \mathbf{H}_\lambda \mathbf{f}\|^2 + \sigma^2 \{\text{tr}(\mathbf{H}_\lambda \mathbf{H}_\lambda^T) - 2\text{tr}(\mathbf{H}_\lambda) + n\}. \end{aligned}$$

Use this expression, I can construct an unbiased estimator for $\text{risk}(\lambda)$ (when σ^2 is known) in the following

$$\widehat{\text{risk}}(\lambda) = \|\mathbf{y} - \hat{\mathbf{f}}_\lambda\|^2 + \sigma^2 \{2\text{tr}(\mathbf{H}_\lambda) - n\}.$$

To make our algorithm more efficient, I did some tricks to speed up the computation. When searching for the optimum λ , I have to fit the model repeatedly for different λ and compute the fitted values through

$$\hat{\mathbf{f}}_\lambda = X(X^T X + \lambda D)^{-1} X^T \mathbf{y}$$

where X is design matrix $x = \begin{pmatrix} 1 & x_1 & x_1^2 & \dots & x_1^p & (x_1 - t_1)_+^p & \dots & (x_1 - t_k)_+^p \\ \vdots & \vdots & \vdots & \dots & \vdots & \vdots & & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^p & (x_n - t_1)_+^p & \dots & (x_n - t_k)_+^p \end{pmatrix}$

t_1, t_2, \dots, t_k are knots

D is a penalty matrix = $\text{diag}\{0, \dots, 0, 1, \dots, 1\}$

the first p+1 are zero and the last k elements are 1.

So I need to compute the inverse matrix for each λ which slows our algorithm down. To keep us from this time consuming step, the cholesky decomposition and eigen decomposition are used. Let the cholesky decomposition of $X^T X$ be $C^T C$, where C is a upper triangular square matrix. Since I use the equally spread knots¹, X is full rank matrix and hence C is invertible. Also, let $C^{-T} = (C^{-1})^T$, then

$$\begin{aligned} X(X^T X + \lambda D)^{-1} X^T &= X(C^T C + \lambda D)^{-1} X^T \\ &= X C^{-1} (I + \lambda C^{-T} D C^{-1})^{-1} C^{-T} X^T \end{aligned}$$

Further, I use the eigen decomposition of $C^{-T} D C^{-1}$, which is $C^{-T} D C^{-1} = P Q P^T$, where Q is a diagonal matrix and P is an orthogonal matrix. Then,

$$\begin{aligned} X(X^T X + \lambda D)^{-1} X^T &= X C^{-1} (I + \lambda P Q P^T)^{-1} C^{-T} X^T \\ &= X C^{-1} P (I + \lambda Q)^{-1} P^T C^{-T} X^T \end{aligned}$$

Let $G = X C^{-1} P$ and thus $X(X^T X + \lambda D)^{-1} X^T = G (I + \lambda Q)^{-1} G^T$, where $(I + \lambda Q)$ is a diagonal matrix and inverting it just the same as inverting the diagonal elements, this will keep us from the complicating matrix inversion step. And since G is the same for all λ , I only need to calculate once. Hence, the computation is much more easier than the original setting.

¹In penalized regression splines, the knots do not necessary to equally spread. However, I place the knots equal-spaced within the domain of the data, ie the knots are placed at the (k/30)th quantiles of the data, for k=1,...,30

(d) The simulation study are based on five criteria to do the fitting. They are cross-validation (CV), generalized cross-validation (GCV), corrected AIC (AIC_c), unbiased risk estimator with unknown sigma (risk1) and unbiased risk estimator with known sigma (risk2). With the simulation setting given in Thomas C. M. Lee (2001), the simulation results are summarized in Figure 1, 2, 3 and 4. We can see that all criteria provide the similar result. The mean $\log_e r$ value is always less than 0.2. All fittings are acceptable even we have so many different simulation settings.

Appendix: Tables and Figures

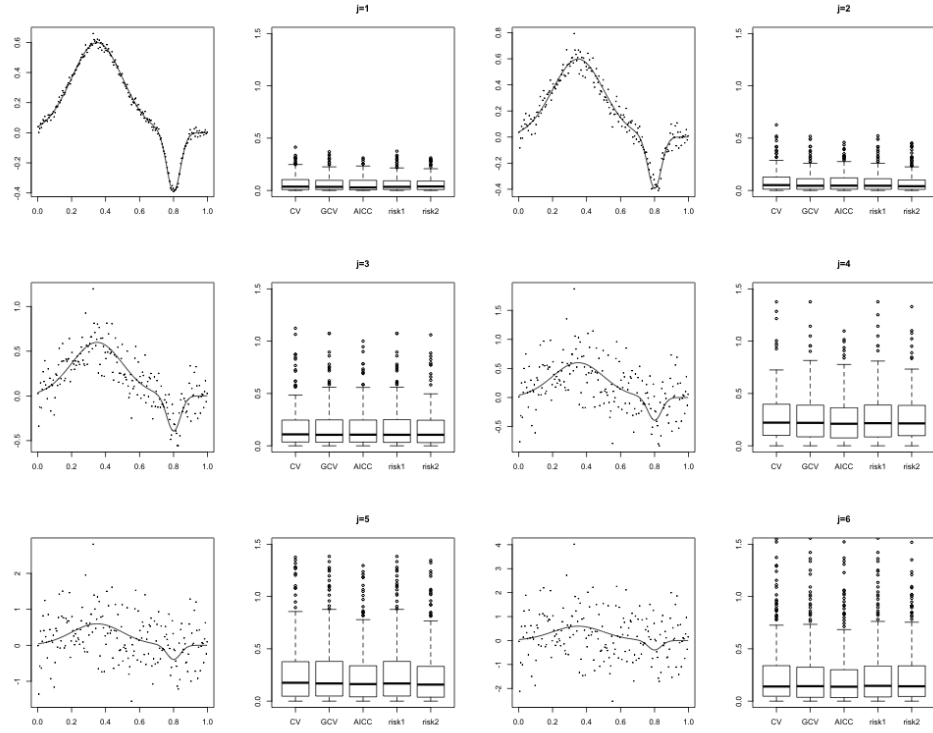


Figure 1: Results for changing the noise level factor. In each pair of panels the left-half displays one typical simulated data set together with the true regression function. The right-half are the boxplots of the $\log_e r$ values for, from left to right, CV, GCV, AIC_c , $risk_1$ and $risk_2$.

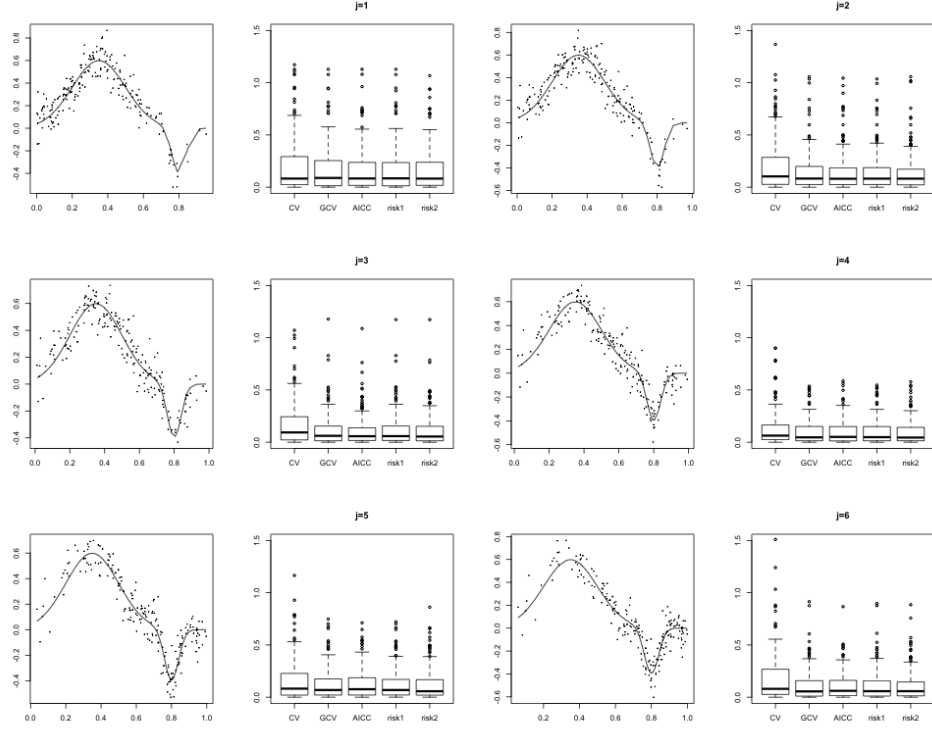


Figure 2: Results for changing the design density level factor. In each pair of panels the left-half displays one typical simulated data set together with the true regression function. The right-half are the boxplots of the $\log_e r$ values for, from left to right, CV, GCV, AIC_c , $risk_1$ and $risk_2$.

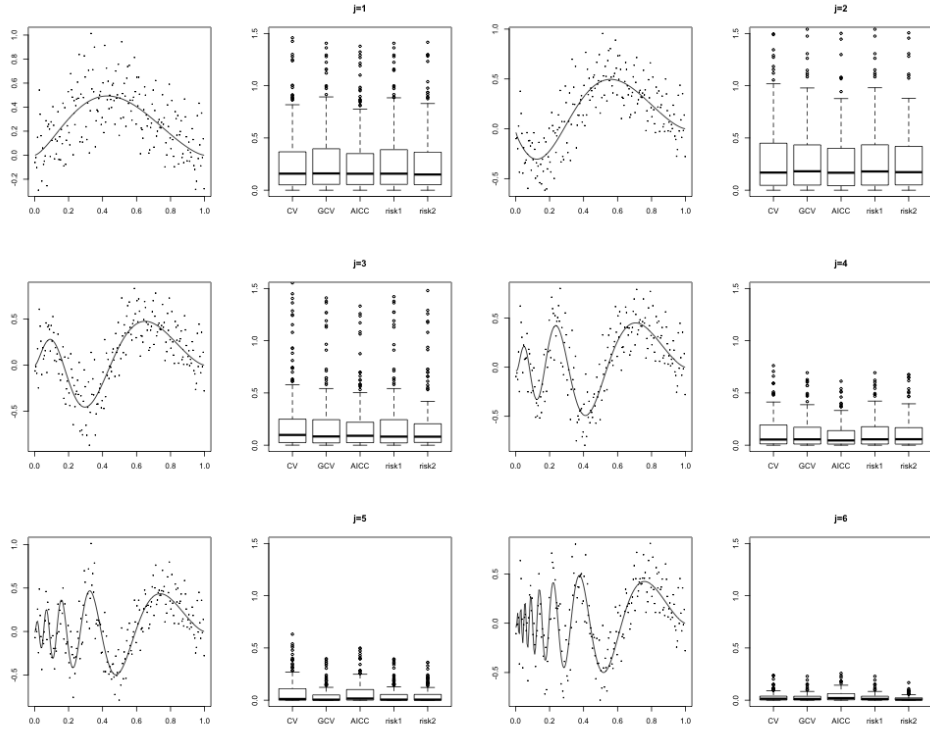


Figure 3: Results for changing the spatial variation level factor. In each pair of panels the left-half displays one typical simulated data set together with the true regression function. The right-half are the boxplots of the $\log_e r$ values for, from left to right, CV, GCV, AIC_c , $risk_1$ and $risk_2$.

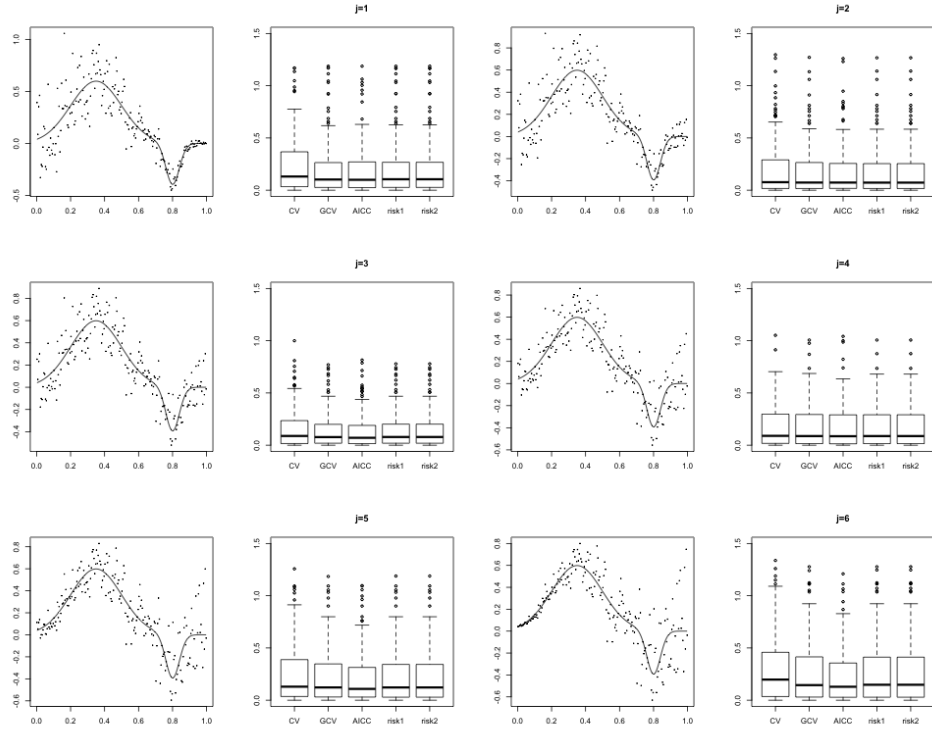


Figure 4: Results for changing the variation function level factor. In each pair of panels the left-half displays one typical simulated data set together with the true regression function. The right-half are the boxplots of the $\log_e r$ values for, from left to right, CV, GCV, AIC_c , $risk_1$ and $risk_2$.