
L^γ Penalty Models

Computation And Applications

Part II

Wenjiang Fu

wfu@stat.tamu.edu

<http://stat.tamu.edu/~wfu>

Department of Statistics, Texas A&M University

SELECTION OF TUNING PARAMETER

Purpose: to achieve best estimation and prediction

Methods: leave-one-out cross-validation (CV), generalized cross-validation (GCV), etc.

Idea: to fit the model well while penalizing on the model size to prevent overfitting.

SELECTION OF TUNING PARAMETER

CV:

Given sample $S = [(x_1, y_1), \dots, (x_n, y_n)]$. Leave one obs. (x_i, y_i) out, and fit model f based on remaining sample $S^{(-i)}$.

Predict y_i with $y_i^* = f_{S^{(-i)}}(x_i)$.

Define $CV = n^{-1} \sum_{i=1}^n (y_i^* - y_i)^2$.

- ◇ CV is computationally expensive.
- ◇ CV is numerically unstable due to outliers, etc.

SELECTION OF TUNING PARAMETER

GCV (Craven and Wahba 1979)

For shrinkage model (Fu 1998, Tibshirani 1996)

$$\text{GCV} = \frac{(y - X\beta)^T (y - X\beta)}{n \left(1 - \frac{\text{tr}(H) - n_0}{n}\right)^2}, \quad (5)$$

where $H = X(X^T X + \lambda W^-)^{-1} X^T$ is a projection matrix, W^- is generalized inverse of $W = \text{diag} \left(2|\hat{\beta}_j|^{2-\gamma}/\gamma\right)$ for $\gamma \geq 1$. Let $n_0 = \#\{\hat{\beta}_j = 0\}$ for lasso only.

Effective number of parameters $p(\lambda, \gamma) = \text{tr}(H) - n_0$.
 $p(0, \gamma) = \text{tr}(X(X^T X)^{-1} X^T) = p$, the number of parameters.
 $p(\infty, \gamma) = 0$ as $\lambda \rightarrow \infty$.

SELECTION OF TUNING PARAMETER

Select $\lambda \geq 0$ for fixed $\gamma \geq 1$

For each fixed $\gamma \geq 1$, compute GCV for each of a sequence of $\lambda \geq 0$ between 0 and a moderate number. Select the value of λ that minimizes GCV.

Select $\lambda \geq 0$ and $\gamma \geq 1$

Compute GCV for each point (λ, γ) on a lattice of $[0, \lambda_0] \times [1, 3]$ with a moderate number λ_0 . Select the values of (λ, γ) that minimize GCV surface.

SELECTION OF TUNING PARAMETER

Problem with GCV

GCV (5) favors lasso even if ridge performs better (Fu 1998).

Reason:

GCV (5) emphasizes linear part by taking $\text{tr}(H)$, performs well for linear estimators, such as ridge.

$\hat{\beta}_{\text{brdg}}$ is nonlinear except for $\gamma = 2$.

For orthonormal X case, lasso is piece-wise linear.

GCV performs poorly in selecting λ for $\gamma \neq 2$.

By Taylor expansion,

$$X\hat{\beta} = H(y)y =$$

$$H(y_0)y_0 + \{H(y_0) + H'(y_0)y_0\}(y - y_0) + o(y - y_0).$$

Thus $\text{tr}(H)$ is a linearization.

Account for nonlinearity

To account nonlinearity, modify GCV (5) through $p(\lambda, \gamma)$. RSS accounts the nonlinearity through the estimator $\hat{\beta}_{\text{brdg}}$. Instead of separating linear part from nonlinear part, we pool them together and consider the overall shrinkage effect through a standard shrinkage rate s .

$$s = \frac{\|\hat{\beta}(\lambda, \gamma)\|_{\gamma}}{\|\hat{\beta}^0\|_{\gamma}},$$

where $\|\cdot\|_{\gamma}$ is the L^{γ} -norm of the shrinkage estimator $\hat{\beta}(\lambda, \gamma)$ or the no-shrinkage estimator $\hat{\beta}^0$ with $\gamma \geq 1$. Apparently $0 \leq s \leq 1$.

Nonlinear GCV

Modify the effective number of parameters

$$p(\lambda, \gamma) = ps$$

where p is the number parameters in the model, s is the standard shrinkage rate.

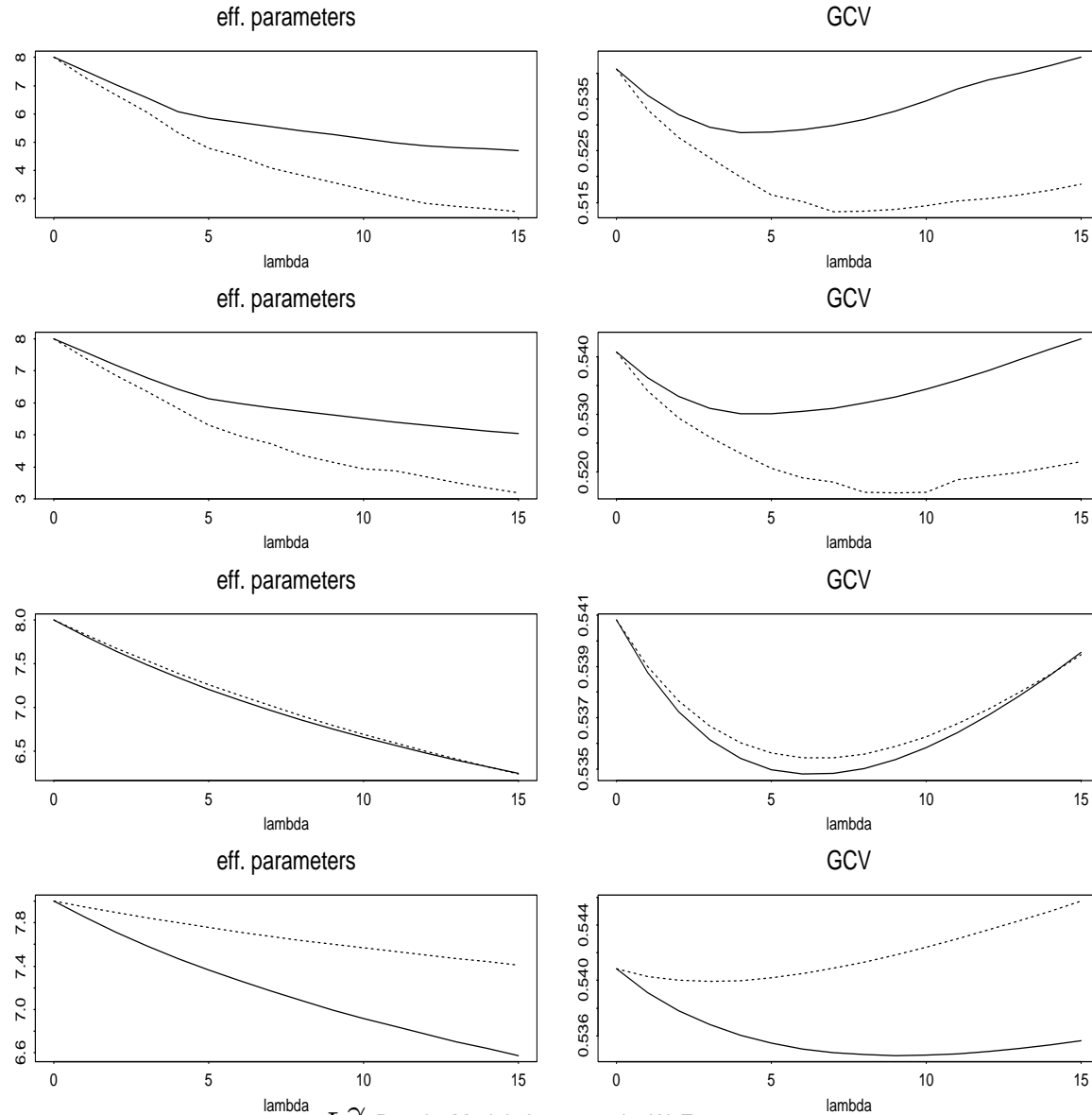
Define the nonlinear GCV as

$$\text{NLGCV} = \frac{\text{RSS}}{n(1 - ps/n)^2}. \quad (6)$$

Refer GCV (5) as linear GCV (LGCV).

SELECTION OF TUNING PARAMETER

Comparison between LGCV and NLGCV for $\gamma = 1, 1.5, 2, 3$.
Solid – NLGCV (6); Dotted – LGCV (5).



SELECTION OF TUNING PARAMETER

Table 2. MSE* in simulation studies with highly collinear X . $n = 10, p = 5$.

model	β^{**}	OLS	LGCV	NLGCV
Lasso	β_1	.1759(.0115)	.1468(.0099)	.0977 (.0118)
	β_2	.0159(.0001)	.0149(.0001)	.0146(.0001)
	β_3	.0618(.0015)	.0534(.0014)	.0389 (.0014)
ridge	β_1	.1679(.0111)	.0898(.0120)	.0821(.0114)
	β_2	.0162(.0001)	.0132(.0001)	.0125(.0001)
	β_3	.0613(.0015)	.0346(.0015)	.0314(.0013)

* $\text{MSE} = (\hat{\beta} - \beta)^T (X^T X) (\hat{\beta} - \beta)$.

** $\beta_1 = (0.5, 1, -0.2, 0, 0)$, $\beta_2 = (1, 0.2, -0.01, -0.5, 0.02)$
and $\beta_3 = (1, 0, 0, 0, 0)$.

SELECTION OF TUNING PARAMETER

Table 3. Comparison of minimum NLGCV by γ
for prostate cancer data

γ	NLGCV*	λ^{**}
1	0.5285	4.33
1.1	0.5300	4.51
2	0.5348	6.36
3	0.5346	9.15

* Value of the minimum NLGCV for fixed γ ;

** Value of λ that minimizes NLGCV for fixed γ .

Conclusion: no γ value dominates the NLGCV. No selection for $\gamma \geq 1$.

Why no selection for γ .

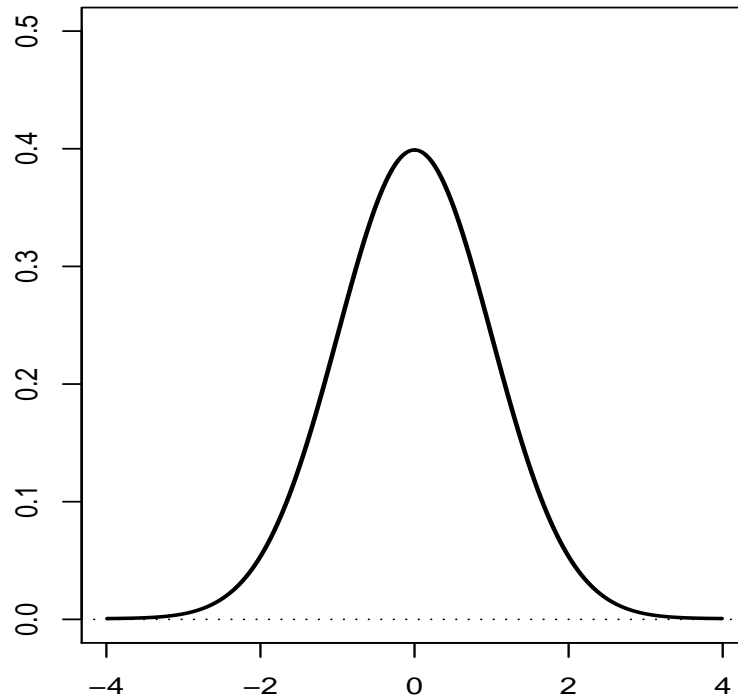
- Bayesian interpretation of L^γ penalty.
 - ◇ $\gamma = 2$: Gaussian prior.
 - ◇ $\gamma = 1$: Laplacian prior.
 - ◇ $\gamma > 1$: complex prior.
- Selecting λ is to select window size for fixed γ .
- Selecting γ is to select prior distribution.
 - ◇ For given data, β may be generated from one prior, say $\gamma = 1.5$.
 - ◇ Prior distributions overlap largely.
 - ◇ Same β may be generated from different priors.
- Conclusion: no selection between priors unless using Bayesian hierarchical model.

SELECTION OF TUNING PARAMETER

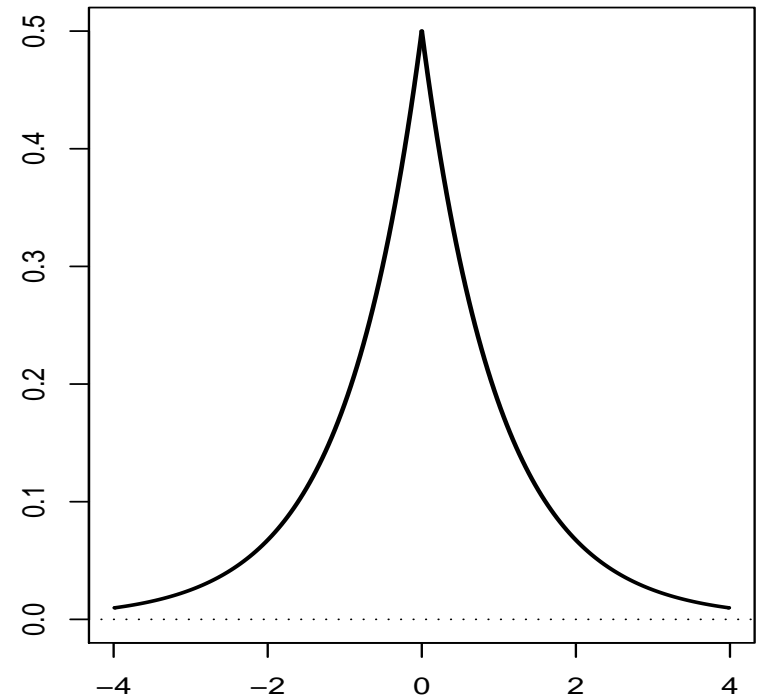
Penalty function as Bayesian prior

$$(\beta|y) \sim C \exp \left\{ -\frac{1}{2} \left(\text{RSS} + \sum \left| \frac{\beta_j}{\lambda^{-1/\gamma}} \right|^\gamma \right) \right\}$$

Gamma = 2



Gamma = 1



Computation of NLGCV

- ◇ Compute $\hat{\beta}_{ols}$ with no penalty.
- ◇ Compute $\hat{\beta}_{brdg}(\lambda, \gamma)$.
- ◇ Compute the ratio of their L^γ norms for s .
- ◇ Compute NLGCV (6).

X not of full rank

- ◇ $\hat{\beta}_{ols}$ is not unique.
- ◇ Compute the limit $\lim_{\lambda \rightarrow 0+} \hat{\beta}_{rdg}(\lambda) = \hat{\beta}_{rdg}(0+)$.
Existence of the limit is guaranteed (Fu 2000).
- ◇ Define standard shrinkage rate s similarly.

Ridge estimator with orthonormal X

For ridge estimator with orthonormal matrix, $X^T X = I$.

$$tr(H) = tr\{X^T (X^T X + \lambda I)^{-1} X\} = p/(1 + \lambda).$$

$$\|\hat{\beta}_{rdg}\|_2 = (1 + \lambda)^{-1} \sqrt{y^T y}, \|\hat{\beta}^0\|_2 = \sqrt{y^T y}.$$

Hence

$$ps = p \frac{\|\hat{\beta}_{rdg}\|_2}{\|\hat{\beta}^0\|_2} = \frac{p}{1 + \lambda} = tr(H).$$

Therefore, LGCV = NLGCV.

Large sample behavior of $\hat{\beta}_{\text{brdg}}$

Finite samples, $\hat{\beta}_{\text{brdg}}$ is biased and performs well in estimation and prediction.

Large samples, is $\hat{\beta}_{\text{brdg}}$ consistent?

Need to study the asymptotics under penalized least squares criterion: to minimize

$$\sum_{i=1}^n (Y_i - x_i^T \phi)^2 + \lambda_n \sum_{j=1}^p |\phi_j|^\gamma.$$

for given λ_n and $\gamma > 0$ fixed.

Regularity conditions

Design $X = (x_1, \dots, x_n)$. x_i are row vectors.

$$C_n = \frac{1}{n} \sum_{i=1}^n x_i x_i^T \rightarrow C,$$

nonnegative definite constant matrix.

$$\frac{1}{n} \max_{1 \leq i \leq n} x_i^T x_i \rightarrow 0.$$

$$\lambda_n/n \rightarrow \lambda_0 \geq 0 \quad (S1)$$

$$\lambda_n/\sqrt{n} \rightarrow \lambda_0 \geq 0 \quad (S2)$$

(S1): λ_n grows fast but not faster than n .

(S2): λ_n grows slowly and not faster than \sqrt{n} .

Limiting distributions

$$Z(\phi) = (\phi - \beta)^T C(\phi - \beta) + \lambda_0 \sum_{j=1}^p |\phi_j|^\gamma.$$

For $\gamma > 1$:

$$V(u) = -2u^T W + u^T C u + \lambda_0 \sum_{j=1}^p u_j \operatorname{sgn}(\beta_j) |\beta_j|^{\gamma-1}.$$

For $\gamma = 1$:

$$V(u) = -2u^T W + u^T C u + \lambda_0 \sum_{j=1}^p [u_j \operatorname{sgn}(\beta_j) I(\beta_j \neq 0) + |u_j| I(\beta_j = 0)].$$

$$W \sim N(0, C\sigma^2).$$

Consistency

Theorem 2. (Knight and Fu 2000)

If C is nonsingular and (S1) is satisfied, then

$$\hat{\beta}_n \rightarrow_p \operatorname{argmin}(Z).$$

So if $\lambda_n = o(n)$, $\hat{\beta}_n$ is consistent.

Theorem 3. (Knight and Fu 2000)

If C is nonsingular and (S2) is satisfied, then

$$\sqrt{n}(\hat{\beta}_n - \beta) \rightarrow_d \operatorname{argmin}(V).$$

Consistency

Theorem 4. (Knight and Fu 2000)

If C is nonsingular and $\lambda_n/n^{\gamma/2} \rightarrow \lambda_0 \geq 0$ for $\gamma < 1$, then

$$\sqrt{n}(\hat{\beta}_n - \beta) \rightarrow_d \operatorname{argmin}(V),$$

where

$$V(u) = -2u^T W + u^T C u + \lambda_0 \sum_{j=1}^p |u_j|^\gamma I(\beta_j = 0)$$

with $W \sim N(0, C\sigma^2)$.

Asymptotic bias

For $\lambda_0 > 0$, asymptotic bias exists for $\gamma \geq 1$.

For example, ridge ($\gamma = 2$),

$$\sqrt{n}(\hat{\beta}_n - \beta) \rightarrow_d C^{-1}(W - \lambda_0\beta) \sim N(-\lambda_0 C^{-1}\beta, \sigma^2 C^{-1}).$$

But for $\gamma < 1$, it is very different. Non-zero β_j can be estimated without asymptotic bias, meanwhile there is a positive mass to shrink $\beta_j = 0$ to 0.

REFERENCES

- Craven, P. and Wahba, G. (1979). Smoothing noisy data with spline functions. *Numerische Mathematik* 31, 377-403.
- Fu, W.J. (1998) Penalized regressions: the Bridge versus the Lasso, *J. Comp. Grap. Statist.* 7: 397-416.
- Fu, W.J. (2000) Ridge estimator in singular design with application to age-period-cohort analysis of disease rates, *Comm. Statist. - Theor. Meth.*, 29:263-278.
- Fu, W.J. (2005) Nonlinear GCV and quasi-GCV for shrinkage models, *Statist. Plann. Infer.*, in press.
- Knight, K. and Fu, W.J. (2000). Asymptotics for Lasso-type estimators, *Ann. Statist.*, 28:1356-1378.
- Tibshirani, R. (1996). Regression shrinkage and selection via the Lasso, *J. Roy. Statist. Soc. B* 58:267-288.