

ridge estimator is biased \leftarrow cf. \rightarrow OLS est'or $\hat{\beta}_0$ is BLUE, but $\text{tr}[\text{cov}(\hat{\beta}_0)] =$
 $(F^T F + \lambda I)^{-1} F^T E(Z) = E(\hat{\gamma}_\lambda) = \gamma - \lambda(F^T F + \lambda I)^{-1} \gamma \Rightarrow \text{Bias}(\hat{\gamma}_\lambda) = -\lambda(F^T F + \lambda I)^{-1} \gamma$
 $\sigma^2 \sum_{j=1}^p \frac{1}{\lambda_j}$
 $(F^T F + \lambda I) - \lambda I = F^T F$ Total of variances: $\text{tr}[\text{cov}(\hat{\gamma}_\lambda)] = \sigma^2 \sum_j \lambda_j (\lambda_j + \lambda)^{-2}$ use (Δ) in LN p. 15
 Mean Square Error: $\text{cov}(\hat{\beta}_\lambda) = E^*(\hat{\beta}_\lambda \hat{\beta}_\lambda^T) = \sigma^2 (F^T F + \lambda I)^{-1} (F^T F) (F^T F + \lambda I)^{-1}$
 $\text{trace} = \text{sum of eigenvalues}$
 $\text{MSE}(\hat{\gamma}_\lambda) = E[(\hat{\gamma}_\lambda - \gamma)(\hat{\gamma}_\lambda - \gamma)^T] = \text{cov}(\hat{\gamma}_\lambda) + \text{Bias}(\hat{\gamma}_\lambda) \text{Bias}(\hat{\gamma}_\lambda)^T \xrightarrow{\text{tr}} \|\text{Bias}(\hat{\gamma}_\lambda)\|^2$
 $\text{Total Mean Square Error: } \text{tr}[\text{MSE}(\hat{\gamma}_\lambda)] = \text{tr}[\text{cov}(\hat{\gamma}_\lambda)] + \lambda^2 \gamma^T (F^T F + \lambda I)^{-2} \gamma$

The total MSE of ridge estimator can be lower than OLS estimator when strong collinearity exists; the price we pay is, of course, the bias.

- why ridge regression can work? \Rightarrow add additional information to remove collinearity. The following conditions, that all lead to ridge estimator, can offer some insights:

Suppose \exists an $n \times p$ matrix V s.t. $V^T F = 0$ and $V^T Z = 0$. Let $W_{n \times p} = \lambda^{1/2} V (V^T V)^{-1/2}$. Then, (1) $W^T W = \lambda I$, (2) $W^T F = 0$, and (3) $W^T Z = 0$. The OLS estimator of the model: $Z = (F+W) \gamma + \epsilon$ is: \rightarrow OLS est'or based on $G: (G^T G)^{-1} G^T Z$

$Z = F\gamma + \epsilon$ \leftarrow cf. \rightarrow $[(F+W)^T (F+W)]^{-1} (F+W)^T Z = (F^T F + \lambda I)^{-1} F^T Z$
 \Rightarrow suitably disturbing F by a small amount to remove strong collinearity

Consider the model $Z = F\gamma + \epsilon$, where $Z_{2n \times 1} = [Z^T \ 0]^T$, $F_{2n \times p} = [F^T \ W^T]^T$. Its OLS estimator is:
 $(F^T F)^{-1} F^T Z = (F^T F + \lambda I)^{-1} F^T Z$
 $F^T F = F^T F + W^T W = F^T F + \lambda I$

pseudo-data to present additional information

\Rightarrow adding additional "cases" to the data set to remove strong collinearity

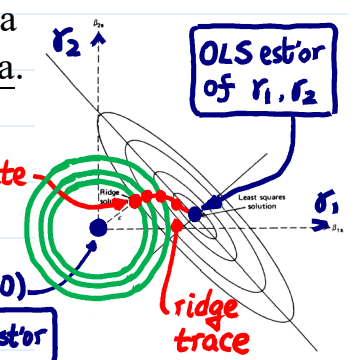
Minimize $RSS = (Z - F\gamma)^T (Z - F\gamma)$ subject to this constraint $\gamma^T \gamma \leq c^2$.
 The solution is the ridge estimator that satisfies $\hat{\gamma}_\lambda^T \hat{\gamma}_\lambda = c^2$.
 This explains why need to standardize X ($X \rightarrow F$)
 Bayesian viewpoint: put a multivariate normal prior $N(0, \lambda^{-1} I)$ on γ .
 Then, the Bayes estimator is the ridge estimator. \Rightarrow choice of a larger λ implies γ were more likely to be small, and vice versa.

L^1 -norm $\sum_{i=1}^p |\hat{\beta}_{i,\lambda}| = c \rightarrow$ LASSO
 double exponential distribution

- an implicit pre-assumption in ridge regression: coefficients (after normalizing) are not likely to be very large
- Reading: Faraway (1st ed.), 9.3 \diamond Further reading: D&S, chapter 17

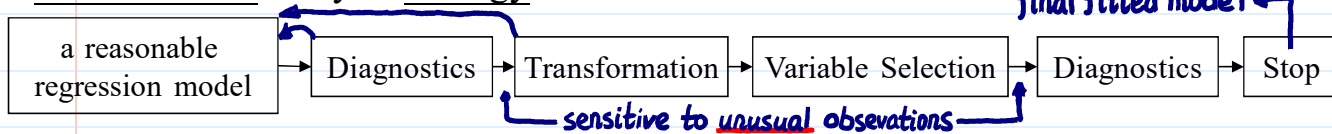
Analysis strategy and model uncertainty

- you have learned
 - Parameter estimation and testing: LS estimator, generalized (weighted) LS, ridge estimator, robust estimator, t -test, F -tests, lack-of-fit, C.I., R^2 , prediction, ...
 - Diagnostics (checking assumptions): such as constant variance, linearity, normality, outliers, influential points, serial correlation, collinearity, ...
 - Transformation: transforming the response and/or the predictors, Box-Cox, polynomial models, broken line, spline, principal component, dummy variables, ...
 - Variable selection: testing-based and criterion-based procedures
- Q: what order should these be done? should procedures be repeated at later stage? when should we stop?



data \rightarrow OLS est'or \rightarrow Bayes est'or

- a recommended analysis strategy:



Note: there is no hard-and-fast rules about how it should be done.

Regression analysis is a search for structure in data. Better to try a variety of orders.

- Danger of doing too much analysis. More transformations, permutations of leaving out influential points and outliers you have done, better fitting model you will find --- however, may lead to over-fitting or no guarantee that the model is a good representation of the underlying system.

➤ avoid complex models for small dataset

➤ try to obtain new data to validate your proposed model

➤ use past experience with similar data to guide the choice of model

CV
 training data
 validation data
 test data

“If you torture the data long enough, it will confess to anything.”

- model multiplicity: Same data can support different models, that sometimes lead to different conclusions. Personal preference, different analysis strategy, or changes in order of analysis components may result in different models. Always try to Bayesian approach take a second independent look at the data. → model averaging for prediction purpose ← cf
- model uncertainty: Usually, inferences are based on the assumption that the selected final model was fixed in advance and so only reflect uncertainty concerning the parameters of that fixed model. **Q**: should we consider the variation caused by model multiplicity? From this viewpoint, the reported standard errors are usually too small.

❖ **Reading:** Faraway (1st ed.), chapter 10