

# Advanced Bayesian Computation

**Rajarshi Guhaniyogi**  
**Winter 2016**

January 15, 2016

# Course Information

- Lectures: MWF 9:30-10:40
- Lecture notes or relevant study materials will be posted every week.
- The course will be graded on two homeworks and one end term project.
- Homework 1: 25%, Homework 2: 25% and End Term: 50%.
- Students taking Satisfactory/Unsatisfactory are required to submit all the homeworks and the final project.
- There will be a 23 minutes presentation for the end term project. I would encourage you to work on the end term project from the late January.
- Lectures will be delivered for 9 weeks. Last week is reserved for the end term presentation.

## High dimensional regression with an emphasis on Bayesian methodology

- Penalized optimization: Ridge regression, lasso, elastic net, adaptive lasso, group lasso.
- Bayesian high dimensional regression:
  - (i) g-prior, two paradoxes, connection with model selection, mixture of g-priors.
  - (ii) Spike and slab prior, detailed discussion, problem with model selection and computation, stochastic search variable selection, issues.
  - (iii) Median probability model in connection with spike and slab prior.
  - (iv) shrinkage estimation, how the name has appeared, motivation, some of the prominent shrinkage priors, Polson and Scott representation.
  - (v) Briefly describe a theoretical result for shrinkage priors.

## **Modeling big data**

- (i) Divide and conquer technique in big data, finding sufficient statistic.
- (ii) Sequential Monte Carlo.
- (iii) Assumed density filtering.
- (iv) Stochastic gradient decent and other applications through stochastic gradient Langevin dynamics.

## **Approximate Bayes method**

- (i) Variational Bayes: Definition, how to compute it.
- (ii) Variational Bayes in nonparametric models.
- (iii) Stochastic variational inference.

# Regression Analysis: An old tool

- Statistical regression is occupying the literature from early 19th century.
- The entire strength of statistics comes from regression analysis.
- With the advancements in computation techniques and various sources of data, regression analysis has been extended to model various situations.
- Our motto is to discuss techniques that makes us up to date with the modern techniques in regression analysis.
- In particular, we will discuss situations where the number of predictors is large.
- Such things typically occur in biomedical applications.

# Linear Regression: Formulation

$$y = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p + \epsilon, \quad \epsilon \sim N(0, \sigma^2)$$

- Different structures of  $\epsilon$  can be accommodated.
- We minimize sum of squared errors to estimate the regression coefficients.

# Understanding Error

- Sum of squared error is a representation of the error in the OLS.
- Sum of squared prediction error is the sum of variance and square of bias.
- Though we only care about the squared prediction error, it becomes helpful to individually understand variance and squared bias.

# Tradeoff between Bias and Variance

- There is a tradeoff between bias and variance in the sense that if model complexity increases, bias decreases, variance increases.
- It is always important to protect from under and over-fitting.
- Important to hit the point with lowest prediction error.



# Gauss Markov Theorem

- Gauss Markov theorem states that among all linear unbiased estimates, OLS has the smallest error.
- There can be some BIASED estimator which is able to provide lower MSE.

# Shrinkage Estimation

- Let OLS estimate is  $\hat{\beta}_j$ . What happens to the MSE if we use an estimator  $\tilde{\beta}_j = \frac{\hat{\beta}_j}{1+\lambda}$ ?

# Shrinkage Estimation

- Let OLS estimate is  $\hat{\beta}_j$ . What happens to the MSE if we use an estimator  $\tilde{\beta}_j = \frac{\hat{\beta}_j}{1+\lambda}$ ?
- **Initially looks like a crazy idea, but lets give it a shot.**
- In particular, can we achieve lower MSE than OLS?
- Yes, we can. But the resulting estimator has to be biased. Whatever we pay for bias is compensated by the variance.
- $\lambda$  that minimizes the error is  $\lambda = \frac{p\sigma^2}{\sum_{j=1}^p \hat{\beta}_j^2}$ .
- Note: As  $\lambda$  becomes big this estimator approaches to 0.

# Shrinkage Estimation

- Charles Stein with his student James found that the estimator  $\beta'_j = \left(1 - \frac{(p-2)\sigma^2}{\sum \hat{\beta}_j^2}\right) \hat{\beta}_j$  has less MSE when  $\sigma^2$  is known.
- Stanley Sclove proposed to shrink the estimator close to zero if we find negative value, i.e.  $\left(1 - \frac{(p-2)\sigma^2}{\sum \hat{\beta}_j^2}\right)^+ \hat{\beta}_j$ .
- If  $\sigma^2$  is unknown, he proposed taking  $\beta'_j = \left(1 - \frac{cRSS}{\sum \hat{\beta}_j^2}\right)^+ \hat{\beta}_j$ , for some constant  $c$ .

# Shrinkage Estimation Contd..

- Note that the F-statistic is given by  $F = \frac{\sum \hat{\beta}_j^2 / p}{RSS / (n-p)}$ .
- Expressing Sclove estimator as  $\beta'_j = \left(1 - \frac{c(n-p)}{pF}\right)^+ \hat{\beta}_j$ , it seems that if the F test statistic is greater than  $c$  then all estimators are set to zero.

# Shrinkage Estimation Contd...

- The above estimation sets all elements to either zero or nonzero.
- Stepwise regression adds or subtracts new variables in the regression if there is an improvement in terms of AIC or BIC.  
 $AIC = n \text{ RSS} + 2 \text{ df}$ ,  $AIC = n \text{ RSS} + \log(n) \text{ df}$ .
- But this is not automated. Is there any method that automates shrinkage?
- What about the shrinkage parameter. Can we use it to estimate stuff?

# Ridge Regression

- In statistical literature, ridge regression was introduced from a completely different perspective.
- Remember, if  $\mathbf{X}$  is the  $n \times p$  matrix and  $\mathbf{y}$  is the  $n \times 1$  responder vector, OLS estimator is given by the solution to the equation  $\mathbf{X}'\mathbf{X}\boldsymbol{\beta} = \mathbf{X}'\mathbf{y}$ .
- Suppose  $\mathbf{X}'\mathbf{X}$  does not have an inverse or the inverse is highly unstable.
- Can happen when  $n < p$  or when columns are highly correlated.
- One idea is to solve  $(\mathbf{X}'\mathbf{X} + \lambda\mathbf{I})\boldsymbol{\beta} = \mathbf{X}'\mathbf{y}$ , with small  $\lambda$ .

# Ridge Regression

- For ridge regression  $\hat{\beta} = (\mathbf{X}'\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}'\mathbf{y}$ .
- Note that  $E(\hat{\beta}) = (\mathbf{X}'\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}'\mathbf{X}\beta \neq \beta$ .
- $Var(\hat{\beta}) = \sigma^2(\mathbf{X}'\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X} + \lambda\mathbf{I})^{-1}$ .
- $\lambda$  is the key parameter. How to choose  $\lambda$ ?



# Generalized Cross Validation to Choose $\lambda$

- **k fold:**

- (i) Divide the data into ten (equal) parts,  $\mathcal{S}_1, \dots, \mathcal{S}_k$ .
- (ii) Set  $\lambda$  on a grid, say  $\lambda \in \{\lambda_1, \dots, \lambda_s\}$ .
- (iii) For every  $\lambda_j$ , use  $\mathcal{S}_{-i_1}$  to fit the model and  $\mathcal{S}_{i_1}$  to calculate model fitting error for  $i_1 = 1, \dots, 10$ .
- (iv) Find the average mean squared error.
- (v) Choose that  $\lambda_j$  which minimizes this error.
- (vi) In general,  $k = 10$  is used.

- **leave one out:**

- (i) When  $n$  is small, generally leave one out cross validation is preferred over the  $k$  fold.
- (ii) Fit the model with  $n - 1$  data points and validate with the  $n$ th one.
- (iii) Repeat it for all sample points to calculate the mean squared error.
- (iv) Choose  $\lambda_j$  that minimizes the error.

# More on Ridge Regression

- Ridge regression will ensure that the coefficients decrease in size.
- In Ridge regression, one does not penalize the intercept as it is in the same scale as the predictors.
- Also predictors can be of vastly different scales. To ensure fair shrinkage to all, generally predictors are standardized.
- This also sets the intercept to zero.
- R code to compute ridge regression is attached.

# Variable Selection

- Variable selection means to select important variables which are affecting the response under the regression model.
- For example, there may be a subset of coefficients which are identically zero. The corresponding predictors have no effect on the regression.
- For ridge regression the coefficients are zero only when  $\lambda = \infty$ .
- Therefore ridge regression **can't select variables**.
- It is useful when a lot of coefficients are close to zero.
- It also does not perform well when a lot of coefficients are moderately large.
- Some post-processing steps may be taken to select variables. But is there any model based straightforward way to select variables?

- Lasso is an acronym for least absolute selection and shrinkage operator.
- It combines the good features of ridge regression with variable selection.
- It is competitive in terms of prediction error w.r.t ridge regression.
- Note that the formulation of ridge regression is

$$\arg \min_{\beta} \sum_{i=1}^n (y_i - \mathbf{x}'_i \beta)^2 + \lambda \sum_{j=1}^p |\beta_j|^2$$

Lasso replaces  $l_2$  penalty by the  $l_1$  penalty, i.e.

$$\arg \min_{\beta} \sum_{i=1}^n (y_i - \mathbf{x}'_i \beta)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

- As  $\lambda$  increases less variables are included, might have higher prediction error after certain  $\lambda$ .
- The idea is to choose  $\lambda$  so as to have proper model fit as well as variable selection.
- $\lambda$  is again chosen using generalized cross validation.
- Code for lasso.
- Great thing about lasso is its property of variable selection. Why it happens to lasso and not to ridge?

# Insight into the Geometry of Lasso and Ridge

- the ridge and lasso optimization can be written as the minimization over  $\beta$

$$\|y - X\beta\|^2 \quad \text{subject to} \quad \|\beta\|_2^2 \leq \lambda$$

$$\|y - X\beta\|^2 \quad \text{subject to} \quad \|\beta\|_1 \leq \lambda.$$

The above is equivalent to the optimization problems

$$(\beta - \hat{\beta}_{OLS})' X' X (\beta - \hat{\beta}_{OLS}) \quad \text{subject to} \quad \|\beta\|_2^2 \leq \lambda$$

$$(\beta - \hat{\beta}_{OLS})' X' X (\beta - \hat{\beta}_{OLS}) \quad \text{subject to} \quad \|\beta\|_1 \leq \lambda.$$

- OLS corresponds to the unconstrained optimization.
- The shapes of ridge and lasso are discussed in class.

# Elastic net: Motivation

- Variable selection with lasso has two shortcomings.
  - (i) The number of variables selected is bounded by the total number of samples in the dataset.
  - (ii) Lasso fails to perform group variable selection, i.e. if a group of variables are correlated, lasso tends to select only one of them.
- Elastic net is motivated by the above two shortcomings.
- You are throwing a net to catch multiple fishes together.

**Theorem:** Suppose  $x_i = x_j$  and  $J(\beta)$  is a strictly convex function. Suppose  $\hat{\beta}$  is obtained by optimizing the objective function  $\|y - X\beta\|^2 + \lambda J(\beta)$ . Then  $\hat{\beta}_i = \hat{\beta}_j$ .

- Since elastic net penalty is strictly convex, elastic net achieves group variable selection.

- The elastic net forms a hybrid of the  $l_1$  and  $l_2$ .
- The  $l_1$  part of the penalty generates a sparse model.
- The quadratic part of the penalty
  - (i) removes limitation on the number of selected variables;
  - (ii) encourages grouping effect.

$$\mathbf{X}_{(n+p) \times p}^* = \frac{1}{\sqrt{(1 + \lambda_2)}} \begin{pmatrix} \mathbf{X} \\ \sqrt{\lambda_2} \mathbf{I} \end{pmatrix}, \mathbf{y}^* = (\mathbf{y}, \mathbf{0})',$$
$$\gamma = \frac{\lambda_1}{\sqrt{1 + \lambda_2}}, \beta^* = \sqrt{1 + \lambda_2} \beta.$$

The elastic net objective function can be written as

$$\|\mathbf{y}^* - \mathbf{X}^* \beta^*\|^2 + \gamma \|\beta^*\|_1.$$

Thus elastic net can select all  $p$  predictors.



# Connections Between Lasso, Elastic Net and Ridge

- Naive elastic net is given by

$$\hat{\beta}_{elastic} = \arg \min ||\mathbf{y} - \mathbf{X}\beta||^2 + \lambda_1 ||\beta||^2 + \lambda_2 ||\beta||_1$$

- Elastic net penalty can be viewed as  $\sum_{j=1}^p [(1 - \alpha)|\beta_j| + \alpha|\beta_j|^2]$ .
- $\alpha = 0$  gives lasso,  $\alpha = 1$  gives ridge.
- Solution to the above elastic net penalty is known as the naive elastic net. Unfortunately it does not perform well in practice.
- The intuitive reason being double penalization.
- Actual elastic net is scaled naive elastic net estimates,  $\beta(enet) = (1 + \lambda_2)\beta(naive\ enet)$ .

# Adaptive Lasso

- The *adaptive lasso* uses a weighted penalty of the form  $\sum_{j=1}^p w_j |\beta_j|$  where  $w_j = 1/|\hat{\beta}_j|^\nu$ ,  $\hat{\beta}_j$  is the ordinary least squares estimate and  $\nu > 0$ .
- The adaptive lasso yields consistent estimates of the parameters while retaining the attractive properties of lasso. Idea is to favor predictors with univariate strength, to avoid spurious selection of noise predictors.
- When  $p > n$ , can use univariate regression coefficients in place of full least squares estimates.
- In general, when the predictors are correlated it is a good practice to use univariate regression coefficients.
- Adaptive lasso recovers the correct model under milder condition than lasso.
- Computationally it does not add any extra significant burden to lasso computation.

# Group Lasso

- In some problems, the predictors belong to pre-defined groups.
- In this situation it may be desirable to shrink and select the members of a group together. The *group lasso* is one way to achieve this.
- Suppose  $p$  predictors are divided into  $m$  groups, with  $p_j$  number of predictors in group  $j$ ,  $j = 1, \dots, m$ ;  
 $p_1 + \dots + p_m = p$ .
- $\mathbf{X}_j$  matrix corresponding to the  $j$ th group of predictors.
- $\beta_j$  is the vector coefficient corresponding to  $\mathbf{X}_j$ .
- Group lasso minimizes

$$\arg \min_{\beta \in \mathbb{R}^p} \left[ \|\mathbf{y} - \beta_0 \mathbf{1} - \sum_{j=1}^m \mathbf{X}_j \beta_j\|^2 + \lambda \sum_{j=1}^m \sqrt{p_j} \|\beta_j\|_2 \right]$$

# Clustering in High Dimensions: Nonnegative Matrix Factorization

- Given a matrix  $\mathbf{M}_{p \times n}$  and a desired rank  $k \ll \min(n, p)$ , find  $\mathbf{W}_{p \times k}$  and  $\mathbf{H}_{k \times n}$  s.t.  $\mathbf{M} \approx \mathbf{WH}$  by solving an optimization problem  $\min_{\mathbf{W} > 0, \mathbf{H} > 0} \|\mathbf{M} - \mathbf{WH}\|^2$ .
- Why do this when SVD does a better job in approximating  $\mathbf{M}$ .
- If  $\mathbf{M} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}$ , then  $\|\mathbf{M} - \mathbf{U}_k\mathbf{\Sigma}_k\mathbf{V}_k\| \leq \|\mathbf{M} - \mathbf{WH}\|$ .
- Reason to do NMF: For nonnegative data NMF approximation provides better interpretation.

# NMF and K-Means Clustering

- k-means clustering can be written as  $\|M - WH\|^2$ .
- Columns of  $H$  gives us the cluster membership indicators.
- Look at the largest element in each column of  $H$ .
- That sample is included in the corresponding cluster.
- Sometimes to make it similar to the K-means, sparse NMF is employed.

# Penalized Optimization: Unsatisfactory in Predictive Inference

- Penalized optimization is unable to provide predictive inference. Only provides point prediction.
- Typical focus in many scientific applications is uncertainty characterization.
- Different choices of tuning parameters may affect inference considerably.

- If loss function corresponds to a likelihood & penalty to the log prior (up to normalizing constants), then estimates correspond to mode of a Bayesian posterior (MAP estimates).
- Consider the linear regression model with known  $\sigma^2$  and with prior

$$y_i \sim N(\mathbf{x}'_i \boldsymbol{\beta}, \sigma^2), \quad \boldsymbol{\beta}_j \sim \pi_{\boldsymbol{\beta}}.$$

- The log posterior of  $\boldsymbol{\beta}$  upto a constant is

$$-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \sum_{j=1}^p \log(\pi_{\boldsymbol{\beta}}(\boldsymbol{\beta}_j))$$

- Although such estimators correspond to the mode of a Bayesian posterior, they are typically not viewed as Bayesian.
- Bayes estimators  $\hat{\beta}_{\text{Bayes}}$  are defined as the value that minimizes the Bayes risk.
- Bayes risk is the expectation of a loss  $L(\hat{\beta}, \beta)$  averaged over the posterior of  $\beta$ .
- For example, if we choose squared error loss,  $\hat{\beta}$  is the posterior mean.
- MAP is not a Bayes estimator for a reasonable choice of loss function.
- Also, we would like to utilize the whole posterior instead of just using a point estimate.