

Regularization Methods for Dimension Reduction

JEFF GILL

Distinguished Professor

Department of Government, Department of Mathematics & Statistics

Member, Center for Neuroscience and Behavior

Founding Director, Center for Data Science

American University

Big Data Background

- ▶ The era of big data has brought about the proliferation of high-dimensional datasets in various fields, including political science.
- ▶ This wealth of data has the potential to enhance our understanding of complex phenomena and contribute to the development of more accurate and reliable predictive models.
- ▶ However, the vast number of variables involved in these high-dimensional datasets presents a unique set of challenges, including multicollinearity and overfitting.
- ▶ Regularization techniques, such as LASSO and Elastic Net, have become popular in addressing these issues by introducing penalties on regression coefficients, encouraging sparsity, and reducing model complexity.

Social Science Perspectives

- ▶ In political science, sociology, public policy, etc., the focus is often on building models that not only provide accurate predictions but also account for theoretical foundations and causal relationships.
- ▶ This requires striking a delicate balance between maintaining the integrity of theoretically significant variables and preventing overfitting or underfitting in the models.
- ▶ Traditional regularization techniques do not always sufficiently cater to this balance, as they may indiscriminately shrink important variables to zero or near-zero values, potentially undermining the theoretical basis of the models.

Why Penalized Regression

- ▶ What if the sample size (n) is small, but the number of predictors (p) are large?
- ▶ With such a large number of potential predictors, often there are problems with *multicollinearity*.
- ▶ To select a smaller subset of predictors: not only fit as well as the full set of variables, but also contains the more important predictors

- ▶ Topically:

We are almost always interested in accurate prediction and determining which predictors are meaningful

- ▶ This is “big data” in the p sense.
- ▶ Regularization is the simultaneous process of selecting a subset of p and estimating the regression coefficients, done with penalties on the covariates such that the more important ones are featured.

section*Early Penalized Regression

- ▶ Start with the standard Linear regression model:

$$y = \mu \mathbf{1}_n + X\beta + \epsilon$$

...and estimate with *all available* p covariates.

$$p \gg n$$

- n observations on Y and p predictors
- \mathbf{X} : the $n \times p$ matrix of *standardized* regressors
- $\boldsymbol{\epsilon} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$

▷ Consider the Residual Sum of Squares (RSS):

$$\text{RSS} = (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})'(\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})$$

- $\tilde{\mathbf{y}} = \mathbf{y} - \bar{y}\mathbf{1}_n$
- OLS estimators are not unique since \mathbf{X} is not of full rank.
- the variances will be artificially large

Early Penalized Regression

- ▶ Hoerl and Kennard (1970): Ridge regression:
 - ▷ minimizes RSS subject to $\sum_{j=1}^p |\beta_j|^2 \leq t$ (L_2 norm).
 - ▷ improves the prediction performance, but cannot produce a model with only the relevant predictors.

- ▶ Frank and Friedman (1993): Bridge regression:
 - ▷ minimizes RSS subject to $\sum_{j=1}^p |\beta_j|^\gamma \leq t$ with $\gamma \geq 0$
 - ▷ the optimal choice of the parameter γ yields reasonable predictors.

- ▶ Fan and Li (2001): Smoothly Clipped Absolute Deviation (SCAD) penalty:
 - ▷ to reduce bias and to yield continuous solutions
 - ▷ derive asymptotic distribution of the estimator with a fixed tuning parameter
 - ▷ show that the estimator satisfies the oracle property (consistent model selection).

The Basic LASSO

- ▶ Among methods that do both continuous shrinkage and variable selection, a promising technique called the **Least Absolute Shrinkage and Selection Operator** (lasso) was proposed by Tibshirani (1996).
- ▶ The lasso is a penalized least squares procedure that minimizes RSS subject to the non-differentiable constraint expressed in terms of the L_1 norm of the coefficients. That is, the lasso estimator is given by

$$\hat{\boldsymbol{\beta}}_L = \arg \min_{\boldsymbol{\beta}} (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})'(\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta}) + \lambda \sum_{j=1}^p |\beta_j|$$

where $\tilde{\mathbf{y}} = \mathbf{y} - \bar{y}\mathbf{1}_n$, \mathbf{X} is the matrix of standardized regressors and $\lambda \geq 0$ is a tuning parameter.

- ▶ Knight and Fu (2000, Annals) have shown consistency for lasso type estimators with fixed p under some regularity conditions on the design. They obtained the asymptotic normal distribution with a fixed true parameter $\boldsymbol{\beta}$ and local asymptotics, that is, when the true parameter is small but nonzero in finite samples. Also, they derived asymptotic properties of lasso type estimators under nearly singular design matrices.

The Basic LASSO

- ▶ For the computation of the lasso, Osborne *et al.* (2000a) proposed two algorithms:
 - ▷ A compact descent algorithm (solves a constrained optimization problem with constraint linearization) was derived to solve the selection problem for a particular value of the tuning parameter.
 - ▷ A homotopy method (an analytic approximation method for highly nonlinear problems using series expansion) for the tuning parameter was developed to completely describe the possible selection.
- ▶ Efron *et al.* (2004) proposed Least Angle Regression Selection (LARS) for a model selection algorithm.
- ▶ This algorithm is a piecewise linear solution path using a modification of forward stagewise and least angle regression paths.
- ▶ They showed that with a simple modification, the LARS algorithm implements the lasso, and one of the advantages of LARS is the short computation time compared to other methods.

Important Limitations of the classic LASSO

- ▶ If there is high multicollinearity among predictors: ridge regression dominates the lasso in prediction performance.
- ▶ In the $p > n$ case, the lasso cannot select more than n variables.
- ▶ If there is a meaningful ordering of the features: the lasso ignores it.
- ▶ A group of variables or the grouped variable: the lasso tends to select individual variables.
- ▶ For the coefficient shrunk to (near) zero, there are no direct forms of their standard errors.

The Basic LASSO, Contrived Example

```
library(glmnet)
X <- matrix(rnorm(200),10,20)
y <- rnorm(10)
lasso.lm <- glmnet(X,y,family="gaussian",alpha=1)
lasso.lm
```

```
      Df  %Dev  Lambda
1      0   0.00 0.60190
2      1   4.03 0.57450
3      1   7.70 0.54840
4      1  11.04 0.52350
5      1  14.09 0.49970
6      1  16.86 0.47700
7      1  19.39 0.45530
:
93     9 99.83 0.00834
94     9 99.84 0.00796
95     9 99.86 0.00760
96     9 99.87 0.00725
97     9 99.88 0.00692
98     9 99.89 0.00661
99     9 99.90 0.00630
```

- (1) the number of nonzero coefficients (Df), (2) the percent (of null) deviance explained (%Dev), (3) and the value of Lambda.

The Basic LASSO, Contrived Example

```
coef(lasso.lm, s=0.5)
```

```
              s1  
(Intercept) 0.134222  
V1           .  
V2           .  
V3           .  
V4           .  
V5           .  
V6           .  
V7           .  
V8           .  
V9           .  
V10          .  
V11          .  
V12          .  
V13          .  
V14          .  
V15          .  
V16          .  
V17          .  
V18          .  
V19          .  
V20          .
```

The Basic LASSO, Contrived Example

```
coef(lasso.lm, s=0.005)
```

```
              s1  
(Intercept) 0.30546423  
V1           0.31145715  
V2           .  
V3           .  
V4          -0.78251180  
V5           .  
V6           .  
V7           0.40352899  
V8           .  
V9           .  
V10          0.59668660  
V11          -0.20360788  
V12          .  
V13          .  
V14          0.04669743  
V15          -0.35153569  
V16          -0.15161324  
V17          .  
V18          0.15144882  
V19          .  
V20          .
```

The Basic LASSO, State Failures Example

- ▶ These data are collected by the State Failure Task Force (SFTF, Esty *et al.* 1999), which is a U.S. government funded group of interdisciplinary researchers whose objective is to understand and forecast when governments cease to function effectively (usually collapsing in violence and disarray).
- ▶ Through a series of reports they have created a warning system of state failures based on the analysis of a huge collection of covariates (about 1,200) on all independent states around the world with a population of at least 500,000, from 1955 to 1998.
- ▶ Thus the greatest challenge is to consider a vast number of potential model specifications using prior theoretical knowledge and model-fitting comparisons.
- ▶ The final results of the SFTF team are controversial because they end up using only three explanatory variables, democracy, trade openness and infant mortality, to produce a model with about 75% correct predictions of state failure (0/1) using the naïve criteria.
- ▶ Their findings are criticized on substantive grounds for being oversimplified (Millien and Krause 2003, Parris and Kate 2003, Sachs 2001), and on methodological ground for their treatment of missing data and forecasting procedures (King and Zeng 2001).

The Basic LASSO, State Failures Data Setup

- ▶ These data are from the State Failure Task Force, which is a U.S. government funded group of interdisciplinary researchers whose objective is to understand and forecast when governments cease to function effectively (usually collapsing in violence and disarray).
- ▶ Through a series of reports they have created a warning system of state failures based on the analysis of a huge collection of covariates (about 1,200) on all independent states around the world with a population of at least 500,000, from 1955 to 1998.
- ▶ The final results of the SFTF team are controversial because they end up using only three explanatory variables, democracy, trade openness and infant mortality, to produce a model with about 75% correct predictions of state failure (0/1) using the naïve criteria.
- ▶ One consistent criticism of the SFTF approach is the use of all global regions in a single analysis. It is clear to area studies scholars that state failures occur with strong regional explanations that can differ significantly.

The Basic LASSO, State Failures Data Setup

- ▶ State Failures Data for 23 Asian countries:

```
sf.asia <- read.table("ARTICLES/Article.Lasso/Data.State.Failures/asia.dat",  
                    sep=" ", header=TRUE)  
sf.asia <- sf.asia[,-1]  
dim(sf.asia)  
[1] 23 129
```

- ▶ Messy issues with missing data:

```
library(mice)  
m <- 5; covars <- 50  
mice.out <- mice(sf.asia[,1:covars],m)  
mice.array <- array(NA,c(nrow(sf.asia),covars,m))  
for (i in 1:m) mice.array[, ,i] <- as.matrix(complete(mice.out,i))  
for (i in 1:m) { for (j in 1:covars) mice.array[,j,i]  
  <- random.imp.vec( mice.array[,j,i]) }  
sum(is.na(mice.array))  
[1] 0
```

The Basic LASSO, State Failures Data Setup

► Run the LASSO:

```
library(glmnet)
X <- mice.array[,-15,1]
y <- mice.array[,15,1]
y[y > 0] <- 1
lasso.lm <- glmnet(X,y,family=binomial(link=probit),alpha=1)
lasso.lm
```

```
      Df  %Dev  Lambda
1      0   0.00 0.60190
2      1   4.03 0.57450
3      1   7.70 0.54840
4      1  11.04 0.52350
5      1  14.09 0.49970
6      1  16.86 0.47700
7      1  19.39 0.45530
8      1  21.70 0.43460
:
92     9  99.81 0.00873
93     9  99.83 0.00834
94     9  99.84 0.00796
95     9  99.86 0.00760
96     9  99.87 0.00725
97     9  99.88 0.00692
98     9  99.89 0.00661
99     9  99.90 0.00630
```


The Basic LASSO, State Failures Data Setup

```
coef.glmnet(lasso.lm,s=0.05)
```

```
              s1
(Intercept) -0.34062463
V1           0.35651339
V2           .
V3           .
V4           0.27364535
V5           .
V6          -0.04733199
V7          -0.04360697
V8           .
V9          -0.39910277
V10          0.40747002
V11          .
V12          .
V13          .
V14          .
V15          -0.57301993
V16          -0.13147005
V17          .
V18          0.30959403
V19          .
V20          .
```

The Basic LASSO, State Failures Data Setup

```
summary(sf.lm.out)$coefficients
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.1449359	0.4164622	0.3480169	0.7610439
X[, 1]	-1.5169291	0.9988976	-1.5186032	0.2681881
X[, 4]	-1.1655855	1.7075531	-0.6826057	0.5653119
X[, 6]	0.9049659	1.1646306	0.7770412	0.5184506
X[, 7]	-0.1233321	0.9126464	-0.1351367	0.9048772
X[, 15]	-1.1244055	1.4816114	-0.7589072	0.5271529
X[, 16]	-0.6992929	0.6822413	-1.0249935	0.4131493
X[, 18]	-0.1458726	0.6178490	-0.2360975	0.8353328

```
summary(sf.lm.out)$r.squared
```

```
[1] 0.8151111
```

```
summary(sf.lm.out)$fstatistic
```

value	numdf	dendf
1.259615	7.000000	2.000000

Fused LASSO

- ▶ If there exists multicollinearity among predictors, ridge regression dominates the lasso in prediction performance.
- ▶ Also, in the $p > n$ case, the lasso cannot select more than n variables because it is the solution to a convex optimization problem.
- ▶ With meaningful ordering of the features (specification of consecutive predictors), the lasso ignores it.
- ▶ If there is a group of variables among which the pairwise correlations are very high and if we consider the problem of selecting grouped variables for accurate prediction, the lasso tends to select individual variables from the group or the grouped variables (for example, dummy variables).
- ▶ To compensate these limitations of the lasso, Tibshirani *et al.* (2005) introduced the fused lasso.
- ▶ The fused lasso penalizes the L_1 -norm of both the coefficients and their differences:

$$\hat{\boldsymbol{\beta}}_F = \arg \min_{\boldsymbol{\beta}} (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})'(\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta}) + \lambda_1 \sum_{j=1}^p |\beta_j| + \lambda_2 \sum_{j=2}^p |\beta_j - \beta_{j-1}|,$$

where λ_1 and λ_2 are tuning parameters.

Group LASSO

- ▶ For grouped variables, Yuan and Lin (2006) proposed a generalized lasso that is called the group lasso.
- ▶ The group lasso estimator is defined as

$$\hat{\boldsymbol{\beta}}_G = \arg \min_{\boldsymbol{\beta}} \left(\tilde{\mathbf{y}} - \sum_{k=1}^K \mathbf{X}_k \boldsymbol{\beta}_k \right)' \left(\tilde{\mathbf{y}} - \sum_{k=1}^K \mathbf{X}_k \boldsymbol{\beta}_k \right) + \lambda \sum_{k=1}^K \|\boldsymbol{\beta}_k\|_{G_k},$$

where K is the number of groups, $\boldsymbol{\beta}_k$ is the vector of β s in group k , the G_k 's are given positive definite matrices and $\|\boldsymbol{\beta}\|_G = (\boldsymbol{\beta}' G \boldsymbol{\beta})^{1/2}$.

- ▶ In general, $G_k = I_{m_k}$, where m_k is the size of the coefficient vector in group k .
- ▶ This penalty function is intermediate between the L_1 penalty and the L_2 penalty.
- ▶ Yuan and Lin argued that it does variable selection at the group level and is invariant under orthogonal transformations.

Elastic Net

- ▶ Zou and Hastie (2005) proposed the elastic net, a new regularization of the lasso, for an unknown group of variables and for multicollinear predictors.

- ▶ The elastic net estimator can be expressed as

$$\hat{\boldsymbol{\beta}}_{\text{EN}} = \arg \min_{\boldsymbol{\beta}} (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})' (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta}) + \lambda_1 \sum_{j=1}^p |\beta_j| + \lambda_2 \sum_{j=1}^p |\beta_j|^2,$$

where λ_1 and λ_2 are tuning parameters.

- ▶ The elastic net estimator can be interpreted as a stabilized version of the lasso.
- ▶ Thus, it enjoys a sparsity of representation and encourages a grouping effect.
- ▶ Also, it is useful when $p \gg n$.
- ▶ They provided the algorithm LARS-EN to solve the elastic net efficiently based on LARS of Efron *et al.* (2004).

Elastic Net with the State Failures Data, Equal Weighting of Penalties

```
elnet.lm <- glmnet(X,y,family=binomial(link=probit),alpha=0.5)
elnet.lm
```

```
      Df  %Dev Lambda
1      0  0.00 1.04600
2      1  2.52 0.99880
3      1  5.00 0.95340
4      1  7.39 0.91010
5      1  9.73 0.86870
6      1 12.00 0.82920
7      1 14.20 0.79150
8      1 16.34 0.75550
9      1 18.41 0.72120
10     2 20.58 0.68840
:
91    16 96.98 0.01590
92    16 97.13 0.01518
93    16 97.28 0.01449
94    16 97.41 0.01383
95    16 97.54 0.01320
96    16 97.66 0.01260
97    16 97.78 0.01203
98    17 97.88 0.01148
99    17 97.99 0.01096
100   17 98.09 0.01046
```

Elastic Net with the State Failures Data, Very Large Lambda

```
coef.glmnet(elnet.lm,s=200)
```

	Df	%Dev	Lambda
V1			.
V2			.
V3			.
V4			.
V5			.
V6			.
V7			.
V8			.
V9			.
V10			.
:			
V40			.
V41			.
V42			.
V43			.
V44			.
V45			.
V46			.
V47			.
V48			.
V49			.

Adaptive LASSO

- ▶ Fan and Li (2001) showed that the lasso can perform automatic variable selection but it produces biased estimates for the larger coefficients.
- ▶ Thus, they argued that the oracle properties (an estimator that is consistent in variable selection is not necessarily consistent in parameter estimation terms of the asymptotic distribution) do not hold for the lasso.
- ▶ To obtain the oracle property, Zhou (2006) introduced the adaptive lasso estimator as

$$\hat{\beta}_{\text{AL}} = \arg \min_{\beta} (\tilde{\mathbf{y}} - \mathbf{X}\beta)'(\tilde{\mathbf{y}} - \mathbf{X}\beta) + \lambda \sum_{j=1}^p \hat{w}_j |\beta_j|,$$

with the weight vector $\hat{\mathbf{w}} = 1/|\hat{\beta}|^\gamma$ where $\hat{\beta}$ is a \sqrt{n} consistent estimator such as $\hat{\beta}(\text{OLS})$ and $\gamma > 0$.

- ▶ The adaptive lasso enjoys the oracle property (“performs as well as if the true underlying model were given in advance”) and it leads to a near-minimax-optimal estimator (has near the least worst range of risk).

Other LASSOs

- ▶ Kim *et al.* (2006) proposed an extension of the group lasso, called blockwise sparse regression (BSR), and studied it for logistic regression models, Poisson regression and the proportional hazards model.
- ▶ Park and Hastie (2007) introduced a path following algorithm for L_1 regularized generalized linear models, and provided computational solutions along the entire regularization path by using the predictor-corrector method of convex optimization.
- ▶ Balakrishnan and Madigan (2007) proposed a data-driven approach, the lasso with attribute partition search (LAPS) algorithm, by combining the fused lasso and the group lasso in particular types of structure classification problems.
- ▶ Meier *et al.* (2008) presented algorithms which are suitable for very high dimensional problems for solving the convex optimization problems, and showed that the group lasso estimator for logistic regression is statistically consistent with a sparse true underlying structure even if $p \gg n$.

Standard Errors of the LASSO

► Samworth (2003)

- ▷ looks at the relation of pointwise asymptotics of consistency of bootstrap estimators and their finite sample behavior
 - ▷ inconsistent bootstrap estimators may in fact perform better.
 - ▷ Hodges-Lehmann: super-efficient estimator (better than asymptotic estimate)
 - ▷ Stein estimators: lower or equal MSE than the OLS estimator.
- Inconsistent Bootstrap
- ▷ the inconsistent bootstrap can only be improved in a very small neighborhood, and the improvements come at the expense of considerably worse performance outside this neighborhood

Standard Errors of the LASSO

- ▶ Beran (1982) shows
 - ▷ for a superefficient estimator, the bootstrap estimates are not consistent if the true parameter is fixed at the point of superefficiency
- ▶ From Knight and Fu (2002) we deduce
 - ▷ The lasso is *superefficient* if $\beta_j = 0$.
- ▶ Pötscher and Leeb (2007) showed that the finite sample distribution of lasso parameters is a mixture of a singular normal distribution and of an absolutely continuous part, which is the sum of two normal densities, each with a truncated tail at the location of the point mass at 0 so this does not give reasonable estimates for the covariance matrix of β .
- ▶ Then Kyung, Gill, Casella, and Ghosh (2010) proved that

the bootstrap estimates of lasso parameters are inconsistent if $\beta_j = 0$.

Why Take a Bayesian Approach to LASSO Regularization?

- ▶ The advantages of the Bayesian version of LASSOs was forcefully argued in Park and Casella (2008) and further developed in Kyung, Gill, Casella, and Ghosh (2010) where a double exponential (Laplace) prior was introduced.
- ▶ Also the penalty term can thus easily be incorporated into the prior structure, whereas classical penalization requires an additional term.
- ▶ Bayesian regularization can simultaneously estimate the penalty parameters in a Gibbs step along with the rest of the model.
- ▶ Posterior standard deviations from MCMC output provide us with valid standard errors and parameter estimates compared to unstable or underperforming standard errors in classical penalization methods.

Background on the Bayesian LASSO

- ▶ Tibshirani (1996) noted that with the L_1 penalty term in the basic form, the lasso estimates could be interpreted as the Bayes posterior mode under independent Laplace (double-exponential) priors for the β_j s.
- ▶ One of the advantages of the Laplace distribution is that it can be expressed as a scale mixture of normal distributions with independent exponentially distributed variances (Andrews and Mallows, 1974).
- ▶ This connection encouraged a few authors to use Laplace priors in a hierarchical Bayesian approach.
- ▶ Figueiredo (2003) used the Laplace prior to obtain sparsity in supervised learning using an EM algorithm.
- ▶ In the Bayesian setting, the Laplace prior suggests the hierarchical representation of the full model.

The Bayesian LASSO of Park and Casella (2008)

- ▶ Estimation is actually easier with Gibbs sampling for the lasso with the Laplace prior in the hierarchical model.
- ▶ Park and Casella (2008) considered a fully Bayesian analysis using a conditional Laplace prior specification of the form:

$$\pi(\boldsymbol{\beta}|\sigma^2) = \prod_{j=1}^p \frac{\lambda}{2\sigma} e^{-\lambda|\beta_j|/\sigma}$$

and the uninformative scale-invariant marginal prior:

$$\pi(\sigma^2) = 1/\sigma^2.$$

- ▶ They pointed out that conditioning on σ^2 is important because it guarantees a unimodal full posterior.
- ▶ Lack of unimodality slows convergence of the Gibbs sampler and makes point estimates less meaningful.
- ▶ Their point estimate recommendation is the posterior median.

Details on the Bayesian Form

- A fully Bayesian analysis using a conditional Laplace prior specification is of the form:

$$\pi(\boldsymbol{\beta}|\sigma^2) = \prod_{j=1}^p \frac{\lambda}{2\sigma} e^{-\lambda|\beta_j|/\sigma}$$

and the uninformative scale-invariant marginal prior:

$$\pi(\sigma^2) = 1/\sigma^2.$$

- In distributional terms:

$$\begin{aligned}\beta_j \mid \tau_j^2, \sigma^2 &\sim \mathcal{N}(0, \sigma^2 \tau_j^2) \\ \tau_j^2 \mid \lambda^2 &\sim \text{Exponential}\left(\frac{\lambda^2}{2}\right), \text{ for } j = 1, \dots, p, \\ \lambda &\sim \text{Half-Cauchy}(0, 1)\end{aligned}$$

- If we integrate τ_j^2 out, we achieve a double-exponential (Laplace) prior:

$$\beta_j \mid \lambda, \sigma \sim \mathcal{L}\left(0, \frac{\sigma}{\lambda}\right), \text{ for } j = 1, \dots, p$$

More Details on the Bayesian Specification

- The Bayesian formulation of the original lasso, as given in Park and Casella (2008), is given by the following hierarchical model.

$$\begin{aligned}
 \mathbf{y} \mid \mu, \mathbf{X}, \boldsymbol{\beta}, \sigma^2 &\sim N_n(\mu \mathbf{1}_n + \mathbf{X} \boldsymbol{\beta}, \sigma^2 \mathbf{I}_n) \\
 \boldsymbol{\beta} \mid \sigma^2, \tau_1^2, \dots, \tau_p^2 &\sim N_p(\mathbf{0}_p, \sigma^2 \mathbf{D}_\tau), \quad \mathbf{D}_\tau = \text{diag}(\tau_1^2, \dots, \tau_p^2) \\
 \tau_1^2, \dots, \tau_p^2 &\sim \prod_{j=1}^p \frac{\lambda^2}{2} e^{-\lambda^2 \tau_j^2 / 2} d\tau_j^2, \quad \tau_1^2, \dots, \tau_p^2 > 0 \\
 \sigma^2 &\sim \pi(\sigma^2) d\sigma^2 \quad \sigma^2 > 0
 \end{aligned}$$

- The parameter μ may be given an independent, flat prior, and posterior propriety can be maintained.
- After integrating out $\tau_1^2, \dots, \tau_p^2$, the conditional prior on $\boldsymbol{\beta}$ has the desired form.
- Any inverted gamma prior for σ^2 would maintain conjugacy, but here we will use the improper prior density $\pi(\sigma^2) = 1/\sigma^2$, with which we also can maintain propriety.
- The resulting prior on $\boldsymbol{\beta}$ is a Laplace distribution with mean 0 and variance $\sigma^2 \lambda^{-2}$ to assure unimodality of the posterior while the unconditional version (without σ^2) does not.

A New Method: the Bayesian Semi-Protected LASSO

- ▶ Again denote the matrix of standardized predictors by X , the outcome variable by Y , the vector of regression coefficients by β , the residual variance by σ^2 , the hyperparameter for the penalty term (L_1 regularization) by λ^2 and the precision by τ^2 .
- ▶ We use n for the number of observations and p for the number of predictors and we also denote the number of protected variables by n_{prot} .
- ▶ The Bayesian LASSO model with protected variables starts with:

$$Y = X\beta + \epsilon,$$

$$\epsilon_i \sim \mathcal{N}(0, \sigma^2), \quad i = 1, \dots, n,$$

$$\beta_j \sim \mathcal{N}(0, \tau_j^2 \sigma^2), \quad j = 1, \dots, p.$$

The Bayesian Semi-Protected LASSO

- ▶ We separate the τ^2 hyperparameters into two groups: protected variables ($\tau_{\text{protected}}^2$) and non-protected variables ($\tau_{\text{non-protected}}^2$).
- ▶ We assign a gamma prior to the protected variables' precision in line and an exponential prior with rate $\lambda^2/2$ to the non-protected precision.
- ▶ Additionally, we assign priors according to:

$$\tau_{\text{non-protected}}^2 \sim \exp(\lambda^2/2),$$

$$\tau_{\text{protected}}^2 \sim \Gamma(1, 1)$$

$$\sigma^2 \sim \mathcal{U}(0.1, 10)$$

$$\lambda^2 \sim \Gamma(\text{shape} = c_0, \text{rate} = d_0)$$

The LASSO in Political Science Research

- ▶ Political scientists typically prioritize reliable prior findings as a theoretical foundation over maximizing predictive accuracy when constructing empirical models of political phenomenon of interest.
- ▶ However, there is increasing interest in predictive modeling in political based on various machine learning tools.
- ▶ This is apparent in published work using random forests, natural language processing, support vector machines, neural networks/deep networks, and more.
- ▶ Yet, in contrast to many other data science fields where there are literally tens of thousands of published articles using the LASSO, there is scant attention to the LASSO for dimension reduction in model specification.
- ▶ In fact we find less than 15 meaningful applications of the LASSO published in political science journals.

The LASSO in Political Science Research

- ▶ So why is the field shy about regularization tools when it appears to fully embrace other machine learning methods?
- ▶ We suspect that it is tied to a traditional adherence to variable selection in settings where the literature has settled on a set of explanations that have reliably endured over time.
- ▶ Given that other machine learning tools are fully embraced by political scientists, there must be something different about regularization tools as a way to reduce the dimension of explanatory variables in model specification.

How Can the Protected Regularization Method Help in Variable Selection?

- ▶ Regularization can sometimes lead to the exclusion of theoretically relevant variables, which may be crucial for understanding the underlying relationships within the data.
- ▶ The protected regularization method offers a solution to this problem by allowing researchers to designate variables that should be protected from shrinkage based on theoretical relevance.
- ▶ This approach is particularly useful when dealing with highly correlated variables that have predictive power, as it enables researchers to prioritize those variables that are most meaningful and relevant to their research questions.
- ▶ A key contribution is a set of measurements for the *cost of protection* in predictive and fit terms.

2020 ANES Application

- ▶ The goal is to predict vote choice in the US presidential election while protecting theoretically important variables from shrinkage.
- ▶ If you're only caring about the prediction, you are not caring about generalization, and vice-versa.
- ▶ Most people using ANES are trying to be in the generalization world, so there's a dichotomy here:
 - ▷ traditional literature using ANES uses theoretically agreed-upon variables, and then they add their own variable; this literature usually aims to establish semi-causal relationships between variables of interest and the outcome variable.
 - ▷ the modern prediction literature has NO preference on variable selection at all and the aim is to predict the vote outcome.
- ▶ We offer a third (general) range of alternatives with our semi-protection.

MSE and BIC Comparison for 3 Models (ANES)

Model	MSE_training	MSE_testing	BIC
No Protection	906.3350	915.3136	79482.38
Partial Protection	910.0663	918.9789	80108.71
Full Protection	922.4238	929.7372	81657.27

Additional Results as of Yesterday

Article	Metric	No-Protection	Partial-Protection	Full-Protection
Simulation data	MSE	94.16159	138.7643	160.0497
politics-social-norms-mask-wearing	MSE	0.7106237	0.7118802	0.7107989
protests-political-attitudes	MSE	0.9560879	0.9568908	0.9572809
Predicting Local Violence	MSE	1.057366	0.9421069	0.9572809
repression-political-loyalty	MSE	0.3697359	0.3620904	0.3590226
politicians-misconduct-colombia	AUC	0.7159553	0.7166278	0.6652021
health-behavior-policy-attitudes-covid	AUC	0.6132714	0.6125319	0.6121527
forced-immigration	MSE	7418180	7467819	7675015
ANES_continuous	MSE	915.3136	918.9789	929.7372

Other Bayesian LASSOs

- ▶ Yuan and Lin (2005) give an empirical Bayes method for variable selection and estimation in linear regression models using approximations to posterior model probabilities that are based on orthogonal designs.
- ▶ Their method is based on a hierarchical Bayesian formulation with Laplace prior and showed that the empirical Bayes estimator is closely related to the lasso estimator.
- ▶ Genkin *et al.* (2007) presented a simple Bayesian logistic lasso with Laplace prior to avoid overfitting and produce sparse predictive models for text data, and
- ▶ Raman *et al.* (2009) used a Bayesian group lasso for a contingency table analysis.

Basic Hierarchical Specification

- Consider hierarchical models of the form

$$\mathbf{y} \mid \mu, \mathbf{X}, \boldsymbol{\beta}, \sigma^2 \sim N_n(\mu \mathbf{1}_n + \mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n), \quad \boldsymbol{\beta} \sim N(0, \Sigma_\beta),$$

where Σ_β is parametrized with τ_i s that are given gamma priors.

- It is important to carefully parameterize Σ_β to obtain the lassos.
- For all the lassos, with the exception of the elastic net, λ and $\boldsymbol{\beta}$ are conditionally independent given the τ_i s, leading to a straightforward Gibbs sampler.

Hierarchical Models and Gibbs Samplers

- A general version of the lasso model can be expressed as

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})'(\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta}) + \lambda_1 h_1(\boldsymbol{\beta}) + \lambda_2 h_2(\boldsymbol{\beta}) \quad \lambda_1, \lambda_2 > 0,$$

where the specific choices of $h_1(\boldsymbol{\beta})$ and $h_2(\boldsymbol{\beta})$ are given by:

Model	λ_1	λ_2	$h_1(\boldsymbol{\beta})$	$h_2(\boldsymbol{\beta})$
lasso	λ	0	$\sum_{j=1}^p \beta_j $	0
Group lasso	λ	0	$\sum_{k=1}^K \ \boldsymbol{\beta}\ _G$	0
			positive definite G_k 's and $\ \boldsymbol{\beta}\ _G = (\boldsymbol{\beta}' G \boldsymbol{\beta})^{1/2}$	
Fused lasso	λ_1	λ_2	$\sum_{j=1}^p \beta_j $	$\sum_{j=2}^p \beta_j - \beta_{j-1} $
Elastic net	λ_1	λ_2	$\sum_{j=1}^p \beta_j $	$\sum_{j=2}^p \beta_j ^2$

Hierarchical Models

- ▶ Now consider how to represent lassos as a conjugate Bayesian hierarchy.
- ▶ For each model we need an unconditional prior on β , and how to represent it as a normal mixture with $\beta \sim N(0, \Sigma_\beta)$, where Σ_β is parametrized with τ_i s.
- ▶ We only need to specify the covariance matrix of β denoted by Σ_β , and the distribution of the τ_i .
- ▶ For the lasso, group lasso, and fused lasso, the covariance matrix Σ_β contains only τ_i s and no λ s. This not only results in β and λ being conditionally independent, it is important for the Gibbs sampler as it results in gamma conditionals for the λ s.
- ▶ This is not the case for the elastic net; to accommodate the squared term we will need to put λ_2 in Σ_β . This can easily be done within the Gibbs sampler. models; details on posterior distributions and Gibbs sampling are left to Appendix ??

Hierarchical Group Lasso

- In penalized linear regression with the group lasso, the conditional prior of $\beta|\sigma^2$ can be expressed as

$$\pi(\beta|\sigma^2) \propto \exp \left(-\frac{\lambda}{\sigma} \sum_{k=1}^K \|\beta_{G_k}\| \right).$$

- This prior can be attained as a gamma mixture of normals, leading to the group lasso hierarchy

$$\begin{aligned} \mathbf{y} \mid \mu, \mathbf{X}, \beta, \sigma^2 &\sim N_n(\mu \mathbf{1}_n + \mathbf{X}\beta, \sigma^2 \mathbf{I}_n) \\ \beta_{G_k} \mid \sigma^2, \tau_k^2 &\stackrel{ind}{\sim} N_{m_k}(\mathbf{0}, \sigma^2 \tau_k^2 \mathbf{I}_{m_k}) \\ \tau_k^2 &\stackrel{ind}{\sim} \text{gamma} \left(\frac{m_k + 1}{2}, \frac{\lambda^2}{2} \right) \quad \text{for } k = 1, \dots, K \end{aligned}$$

where m_k is the dimension of G_k , the grouping matrix.

- For the group lasso we need to use a gamma prior on the τ_i , but the calculations are quite similar to those of the ordinary lasso.

Hierarchical Fused Lasso

- In penalized linear regression with the fused lasso, the conditional prior of $\beta|\sigma^2$ can be expressed as:

$$\pi(\beta|\sigma^2) \propto \exp \left(-\frac{\lambda_1}{\sigma} \sum_{j=1}^p |\beta_j| - \frac{\lambda_2}{\sigma} \sum_{j=1}^{p-1} |\beta_{j+1} - \beta_j| \right).$$

- This prior can also be obtained as a gamma mixture of normals, leading to the hierarchical model:

$$\begin{aligned} \mathbf{y} \mid \mu, \mathbf{X}, \beta, \sigma^2 &\sim N_n(\mu \mathbf{1}_n + \mathbf{X}\beta, \sigma^2 \mathbf{I}_n) \\ \beta \mid \sigma^2, \tau_1^2, \dots, \tau_p^2, \omega_1^2, \dots, \omega_{p-1}^2 &\sim N_p(\mathbf{0}, \sigma^2 \Sigma_\beta) \\ \tau_1^2, \dots, \tau_p^2 &\sim \prod_{j=1}^p \frac{\lambda_1^2}{2} e^{-\lambda_1 \tau_j^2 / 2} d\tau_j^2, \quad \tau_1^2, \dots, \tau_p^2 > 0 \\ \omega_1^2, \dots, \omega_{p-1}^2 &\sim \prod_{j=1}^{p-1} \frac{\lambda_2^2}{2} e^{-\lambda_2 \omega_j^2 / 2} d\omega_j^2, \quad \omega_1^2, \dots, \omega_{p-1}^2 > 0 \end{aligned}$$

Hierarchical Fused Lasso

- Here the $\tau_1^2, \dots, \tau_p^2, \omega_1^2, \dots, \omega_{p-1}^2$ are mutually independent, and Σ_β is a tridiagonal matrix with:

$$\text{Main diagonal} = \left\{ \frac{1}{\tau_i^2} + \frac{1}{\omega_{i-1}^2} + \frac{1}{\omega_i^2}, i = 1, \dots, p \right\},$$

$$\text{Off diagonals} = \left\{ -\frac{1}{\omega_i^2}, i = 1, \dots, p-1 \right\},$$

where, for convenience, we define $(1/\omega_0^2) = (1/\omega_p^2) = 0$.

- Here for the first time we have correlation in the prior for β , adding some difficulty to the calculations.

Hierarchical Elastic Net

- In penalized linear regression with the elastic net, the conditional prior of $\beta|\sigma^2$ can be expressed as:

$$\pi(\beta|\sigma^2) \propto \exp \left\{ -\frac{\lambda_1}{\sigma} \sum_{j=1}^p |\beta_j| - \frac{\lambda_2}{2\sigma^2} \sum_{j=1}^p \beta_j^2 \right\}.$$

- This prior can also be written as a normal mixture of gammas, leading to the hierarchical model

$$\begin{aligned} \mathbf{y} \mid \mu, \mathbf{X}, \beta, \sigma^2 &\sim N_n(\mu \mathbf{1}_n + \mathbf{X}\beta, \sigma^2 \mathbf{I}_n) \\ \beta \mid \sigma^2, \mathbf{D}_\tau^* &\sim N_p(\mathbf{0}_p, \sigma^2 \mathbf{D}_\tau^*), \\ \tau_1^2, \dots, \tau_p^2 &\sim \prod_{j=1}^p \frac{\lambda_1^2}{2} e^{-\lambda_1^2 \tau_j^2 / 2} d\tau_j^2, \quad \tau_1^2, \dots, \tau_p^2 > 0, \end{aligned}$$

where \mathbf{D}_τ^* is a diagonal matrix with diagonal elements $(\tau_i^{-2} + \lambda_2)^{-1}$, $i = 1, \dots, p$.

- In this case, β is not conditionally independent of λ_2 , as it appears in the covariance matrix.

Tuning Parameters

- ▶ The lassos just discussed all have tuning parameters: λ_1 and λ_2 .
- ▶ Park and Casella (2008) suggested some alternatives based on empirical Bayes using marginal maximum likelihood, putting λ_1 or λ_2 into the Gibbs sampler with an appropriate hyperprior. In this paper,
- ▶ We can use the suggested gamma prior for a proper posterior from Park and Casella (2008), and also for comparison, estimate the tuning parameters with marginal maximum likelihood implemented with an EM/Gibbs algorithm (Casella, 2001).
- ▶ The gamma priors on λ^2 are given by:

$$\pi(\lambda^2) = \frac{\delta^r}{\Gamma(r)} (\lambda^2)^{r-1} e^{-\delta\lambda^2}, \quad (r > 0, \delta > 0).$$

State Failures in Asia Again

- ▶ The final results of the SFTF team are controversial because they end up using only three explanatory variables: democracy, trade openness and infant mortality, to produce a model with about 75% correct predictions of state failure (0/1) using the naïve criteria.
- ▶ In the original data for this subset, there are $p = 128$ explanatory variables with $n = 23$ observations.
- ▶ Here we omit four cases that had greater than 70% missing data, making standard missing data tools inoperable.
- ▶ In addition, 18 variables provided no useful information and were dropped as well.
- ▶ So we have $p = 110$ explanatory variables with $n = 19$ observations to apply to the LARS algorithm, the Bayesian lasso, and the Bayesian Elastic Net, all with a probit link function for the dichotomous outcome of state failure.

State Failures in Asia Again

- ▶ Interestingly, the LARS lasso picks three like the stepwise procedure of the State Failure Task Force, but they are three *different* explanatory variables:
 - ▷ `polxcons`: the level of constraints on the political executive, from low to high
 - ▷ `sftpeind`: an indicator for ethnic war, 0 = none, 1 = at least one
 - ▷ `sftpmmmax`, the maximum yearly conflict magnitude scale
- ▶ For a more detailed explanation of these variables see <http://globalpolicy.gmu.edu/pitf>).
- ▶ Recall that the LARS lasso is a variable weighter not a variable selector, where the weights are either zero or one.
- ▶ Thus the LARS lasso zeros-out 107 variables here in favor of the 3 listed above.

State Failures in Asia Again

- ▶ This table provides the top ten variables by absolute posterior median effect from the Bayesian lasso, and also for these variables in the LARS lasso conclusion.
- ▶ The Bayesian Elastic Net produces results that are virtually indistinguishable from the Bayesian lasso for these data.

Variable	Bayesian Lasso Quantiles				LARS Lasso	
	0.05	0.10	0.50	0.90	0.95	Weight
sftpeind	-0.2387	-0.0888	0.3823	1.6498	2.3219	0.1999
sftpmmmax	-0.3282	-0.1686	0.2257	1.2086	1.6969	0.0307
sftpomag	-0.4095	-0.2380	0.1927	1.1228	1.6081	0.0000
sftpcons	-0.4197	-0.2315	0.1846	1.0907	1.5559	0.1750
sftpnum	-0.4430	-0.2407	0.1657	1.0253	1.4062	0.0000
sftpem1	-0.4421	-0.2636	0.1480	1.0140	1.4609	0.0000
dispop1	-1.3756	-0.9410	-0.1213	0.3031	0.5050	0.0000
sftpeth	-0.5091	-0.2951	0.1194	0.9251	1.3498	0.0000
sftgreg2	-0.4616	-0.2811	0.1137	0.9115	1.3047	0.0000
polpacmp	-0.5071	-0.3106	0.1098	0.8568	1.2240	0.0000

State Failures in Asia Again

- ▶ Every coefficient credible interval of the Bayesian lasso covers zero, including the ten given in the table.
- ▶ This indicates a broad lack of traditional statistical reliability despite the three choices of the LARS lasso.
- ▶ Recall that we cannot produce corresponding credible intervals for the LARS lasso, so the credible intervals from the Bayesian lasso are the only available measure of uncertainty.
- ▶ Thus, given the evidence from the Bayesian lasso, we are inclined to believe that the LARS lasso is overly-optimistic with these choices.
- ▶ This example is interesting because of the difficulty in picking from among the many $p \gg n$ possible right-hand-side variables and the suspect stepwise manner taken by the creators of the data producing only `poldemoc`, `pwtopen`, and `sfxinfm`.
- ▶ The LARS lasso picked none of these three but also picked a very parsimonious set of explanations.
- ▶ The Bayesian lasso, as indicated in the table, finds little support for the the stepwise-chosen result.

State Failures in Asia Again

- ▶ In terms of posterior effect size (absolute value of the posterior median), the Bayesian lasso's top four variables contain the three picked by the LARS lasso. The sole exception is:
 - ▷ `sftpmag`: the magnitude of conflict events of all types,
which is closely related to `sftpmmax`.
- ▶ This is reassuring since it implies that the two approaches are focusing on a small core of potentially important explainers.

State Failures in Asia Again

- ▶ While picking the top ten effect sizes is arbitrary (in the table), there is a noticeable drop in magnitude after this era.
- ▶ The Bayesian lasso therefore brings some additional variables to our attention:
 - ▷ `sftpnum`: the number of critical (negative political events).
 - ▷ `sftpem1`: the ethnic war magnitude indicator number 1.
 - ▷ `dispop1`: the population proportion of the largest politically significant communal group seeking autonomy and subject to discrimination.
 - ▷ `sftpeth`: the ethnic wars score.
 - ▷ `sftgreg2`: the subregion used by `sftf...` scores.
 - ▷ `polpacmp`: a 0-10 point indicator with increasing levels of autocratic governmental control.
- ▶ These variables reinforce the themes in the first first four: ethnic groups, ethnic conflict, magnitude of conflict, and the level of executive control of government, without broad consultation, appear to be important determinants of state failure.

Final Notes on Bayesian Lassos

- ▶ Posterior means will never be exactly zero, so a method of selection based on that assumption cannot be achieved.
- ▶ If selection is desired, one strategy is to set equal to zero any coefficient estimate whose confidence (credible) interval includes zero.
- ▶ As the non-Bayesian lasso cannot produce valid standard errors if the true coefficients are zero, it cannot give any confidence assessment of these selections.
- ▶ Having the MCMC output allows us to summarize the posterior in any manner that we choose – although it is typical to use the posterior mean, we could also use the posterior mode *and* a measure of uncertainty..
- ▶ Putting λ into the Gibbs iterations, the estimates of the regression coefficients are not based on a fixed value, but rather are marginalized over all λ , leading to somewhat of a robustness property.