

# "Honey, I shrunk the parameters": A Comparison of Lasso, Ridge, and Elastic Net Methods

Carolina Alvarez, Edoardo Falchi, and Emily Anne Schwab

Research Module in Econometrics and Statistics

January 18, 2022



# Table of Contents

- 1 Introduction
- 2 Statistical Properties of Regularization Methods
  - Ridge
  - Lasso
  - Elastic Net
- 3 Simulation Studies
  - Comparison of Training MSE
  - Selecting an Optimal Tuning Parameter
- 4 Data Application
  - Data Set Overview
  - Coefficient Paths
  - Model Selection
- 5 Concluding Remarks

# Introduction: Drawbacks to OLS

- Consider the following multivariate linear regression, where  $Y = X\beta + \varepsilon$ :
  - $Y = (y_1, \dots, y_n) \in \mathbb{R}^n$  is the response vector
  - $X \in \mathbb{R}^{n \times p}$  is the predictor variable matrix
  - $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$  is the error vector
- When is OLS a sub-optimal model for prediction?
  - Imperfect multicollinearity
    - $\text{cov}(X_i, X_j) \neq 0$
    - $\beta_{OLS}$  still BLUE, but variance and standard errors will increase.
    - Leads to lower t-statistics.
  - High-dimensionality ( $p > n$ )
    - $\text{Rank}(X) < p$  such that OLS does not generate a unique solution.
    - Variance of estimated coefficients becomes infinitely large.

# Regularization Methods

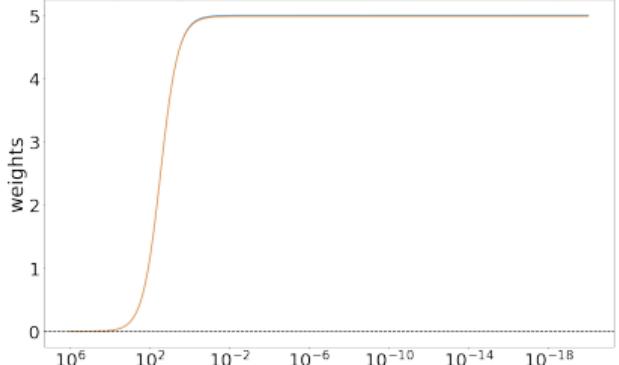
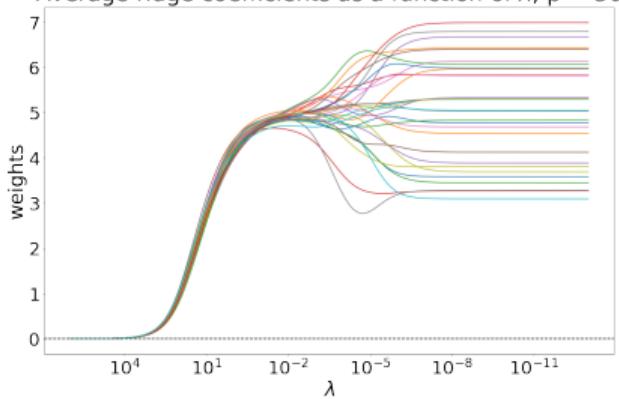
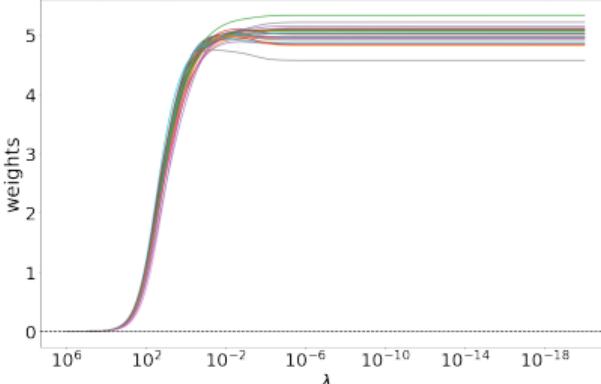
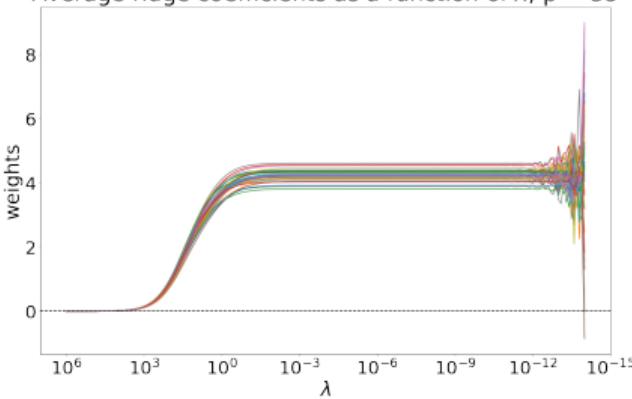
**Solution:** Use regularization method(s) to reduce the variance of the model fit in exchange for a marginal increase in the bias.

# Ridge Regression

Similar to OLS, ridge minimizes the sum of squared residuals (SSR), but includes an  $\ell_2$ -norm shrinkage penalty with a tuning parameter ( $\lambda$ ):

$$\underbrace{\sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2}_{\text{SSR}} + \underbrace{\lambda \sum_{j=1}^p \beta_j^2}_{\text{shrinkage penalty}}. \quad (1)$$

- If  $\lambda = 0$ , we are left with the OLS estimates.
- As  $\lambda$  increases, the estimated ridge coefficients shrink towards zero.

Average ridge coefficients as a function of  $\lambda$ ,  $p = 2$ Average ridge coefficients as a function of  $\lambda$ ,  $p = 30$ Average ridge coefficients as a function of  $\lambda$ ,  $p = 28$ Average ridge coefficients as a function of  $\lambda$ ,  $p = 35$ 

# Ridge Regression

Alternatively,

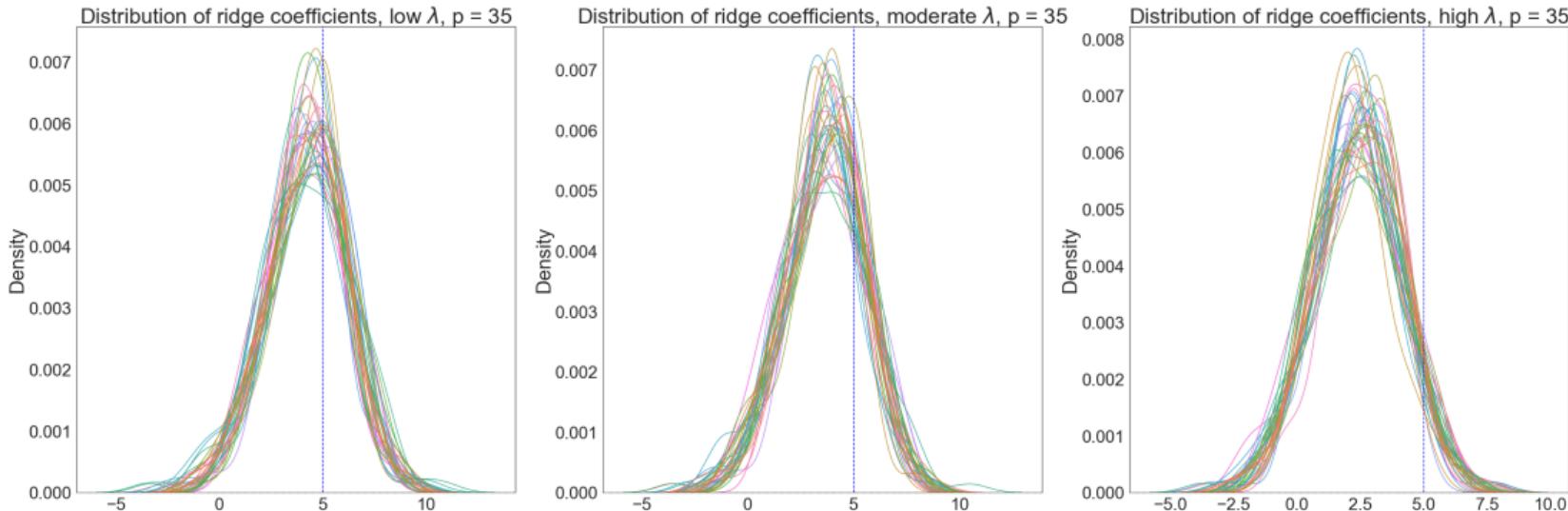
$$\underset{\beta}{\text{minimize}} \left\{ \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} \text{ subject to } \sum_{j=1}^p \beta_j^2 \leq t, \quad (2)$$

where decreasing values of  $t$  indicate an increasingly restrictive optimization constraint.  
 Ridge regression yields a closed-form solution for  $\beta_{ridge}$ :

$$\hat{\beta}_{\lambda}^R = \mathbf{Z}_{\lambda} \hat{\beta} = (\mathbf{X}' \mathbf{X} + \lambda \mathbf{I})^{-1} (\mathbf{X}' \mathbf{X}) \hat{\beta} \quad (3)$$

and for  $\text{var}(\beta_{ridge})$ :

$$\mathbf{V}(\hat{\beta}_{\lambda}^R) = \sigma^2 [\mathbf{X}' \mathbf{X} + \lambda \mathbf{I}]^{-1} (\mathbf{X}' \mathbf{X}) [\mathbf{X}' \mathbf{X} + \lambda \mathbf{I}]^{-1}. \quad (4)$$



# Lasso Regression

Consider the lasso minimization problem, which uses an  $\ell_1$ -norm constraint:

$$\underset{\beta}{\text{minimize}} \left\{ \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j X_{ij} \right)^2 \right\} \text{ subject to } \sum_{j=1}^p |\beta_j| \leq t, \quad (5)$$

- where the  $\hat{\beta}_\lambda^L$  solution is *sparse*: the lasso model holds onto relevant (non-zero) coefficients and sets irrelevant coefficients to **zero**
- Unlike ridge, lasso is a quadratic programming problem that does not have a closed-form solution.

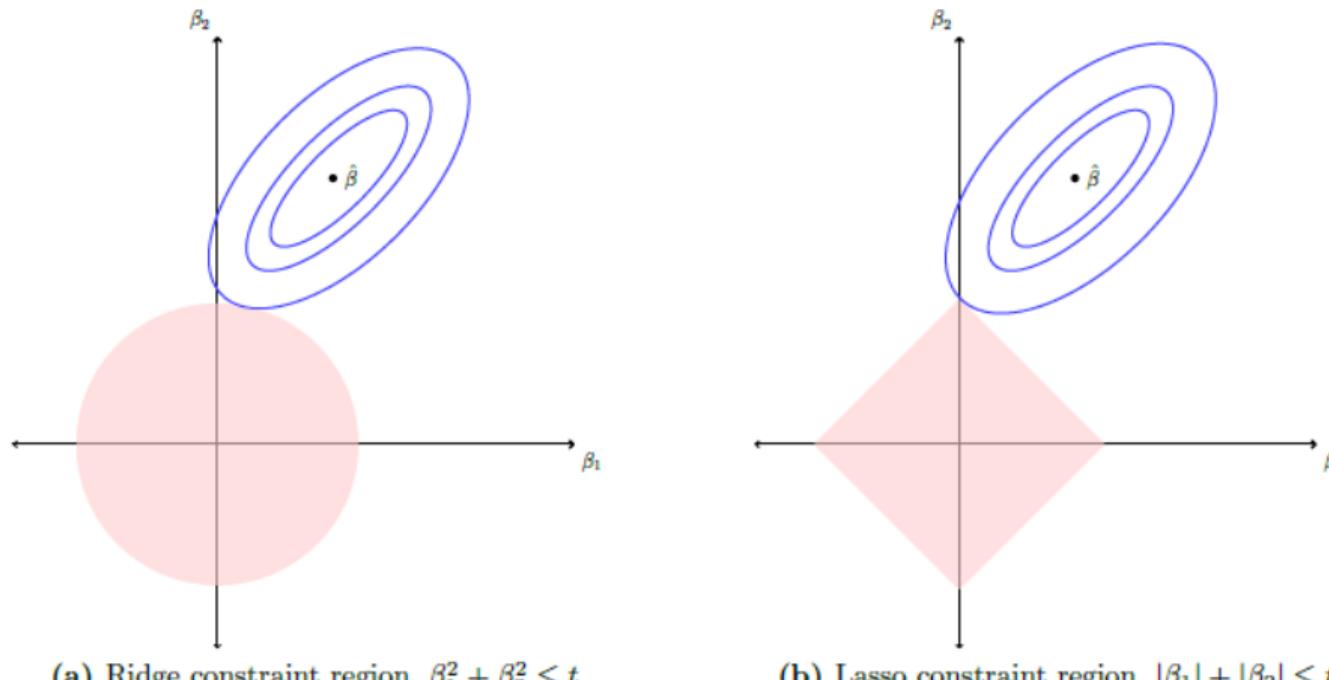
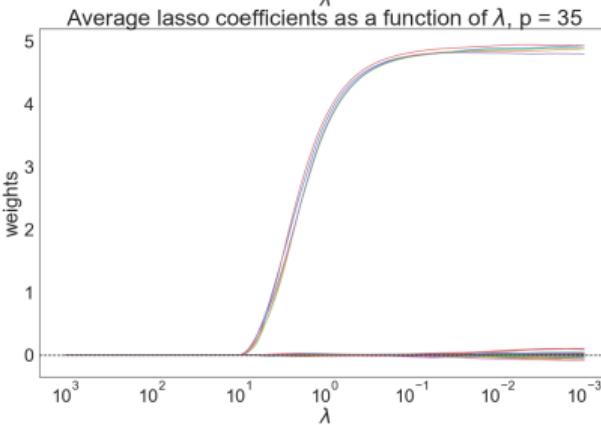
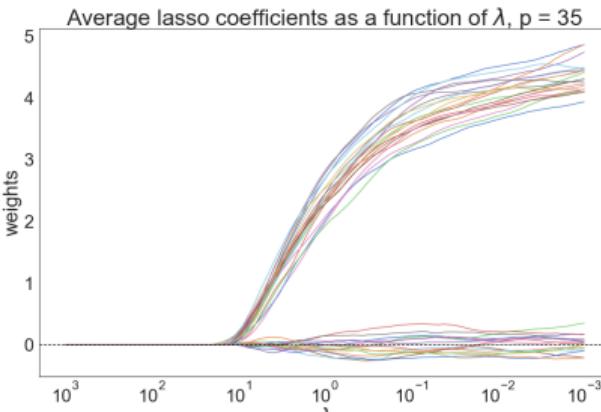
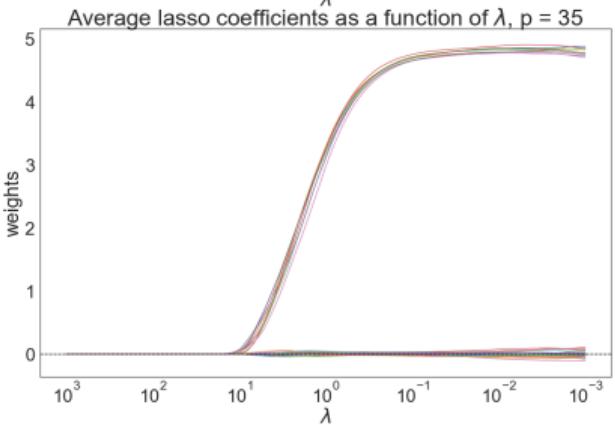
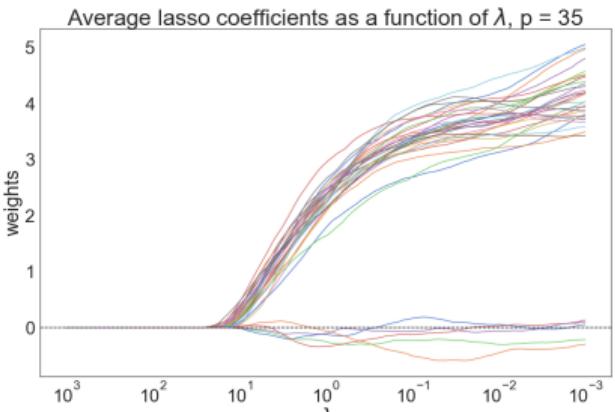


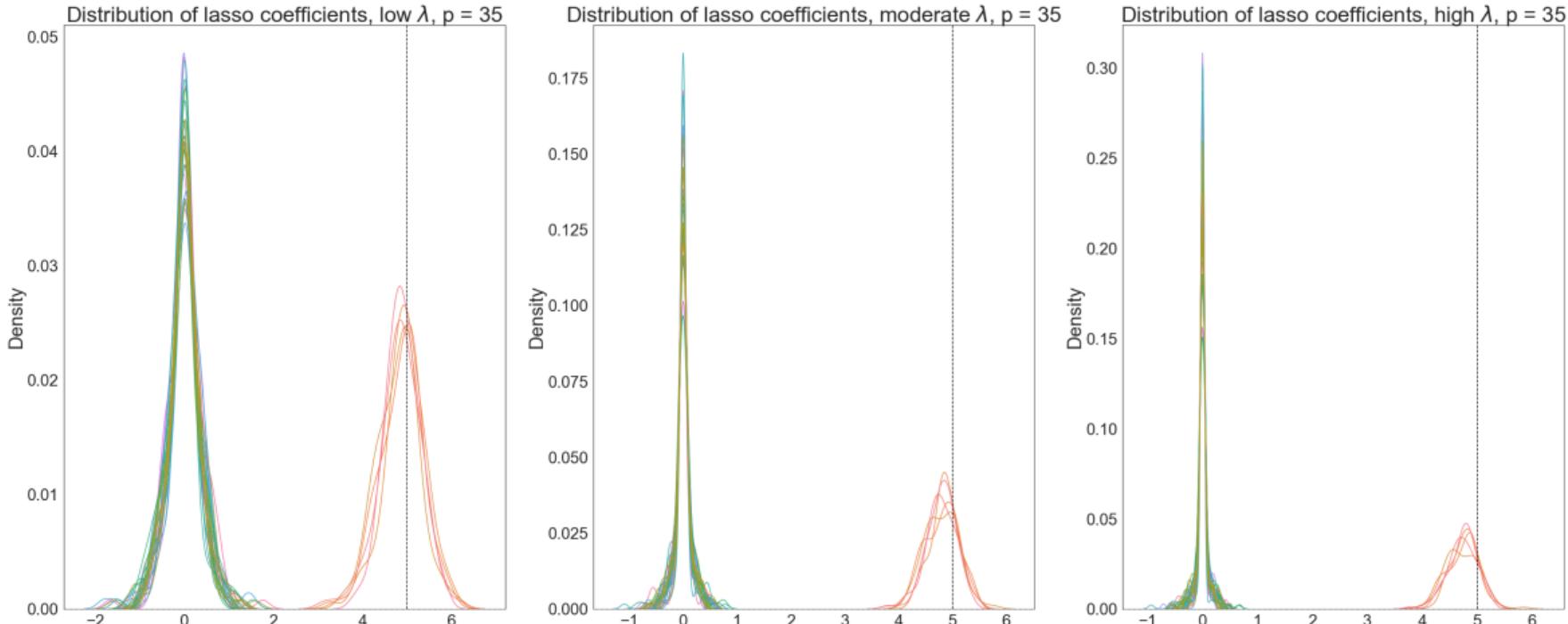
Figure 1: Contours of the SSR (in blue) and the constraint region (in red).

# Lasso Regression

## Limitations:

- Degree of sparsity
  - Lasso does poorly when a large subset of regressors are  $\{j : \beta_j \neq 0\}$  of  $\{1, \dots, p\}$
- With high pairwise correlation among predictors, ridge regression outperforms lasso in terms of prediction.
- Lasso can only select up to  $n$  regressors when  $p > n$ .





# "Naïve" Elastic Net Regression

Consider the "naïve" elastic net minimization problem, which is subject to a constraint that is represented as a convex combination of the ridge and lasso shrinkage penalties:

$$\underset{\beta}{\text{minimize}} \left\{ \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} \text{ subject to } (1 - \alpha) \sum_{j=1}^p |\beta_j| + \alpha \sum_{j=1}^p \beta_j^2 \leq t \quad (6)$$

where  $\alpha$  represents the weight assigned to the ridge penalty.

- Contains the best features of ridge and lasso regression (continuous shrinkage and simultaneous variable selection).
  - For all  $\alpha \in [0, 1]$ , the elastic net penalty function is singular at 0 and strictly convex for all  $\alpha > 0$ .
  - Strict convexity is what sets elastic net apart from lasso.

# "Naïve" Elastic Net Regression

- **Improves over the lasso by**

- overcoming the saturation of variable selection when  $p > n$ . (**variable selection**)
- accounting for highly correlated "grouped" variables. (**continuous shrinkage**)
- strengthening prediction performance if substantial multicollinearity among predictors exists in cases where  $n > p$ .

- **Limitations**

- Double shrinkage does not reduce variance and unnecessarily introduces added bias.
- Hence, works only when its solution is very close to either ridge or lasso.

# Elastic Net Regression

Given the data set  $(y, X)$  and the two fixed non-zero Lagrangian parameters,  $(\lambda_1, \lambda_2)$ , derived from the naïve elastic net's optimization problem, *Lemma 1* of Zou et al. 2005 defines the artificial data set  $(\mathbf{y}^*, \mathbf{X}^*)$ :

$$\hat{\beta}^* = \arg \min_{\hat{\beta}^*} |\mathbf{y}^* - \mathbf{X}^* \beta^*|^2 + \frac{\lambda_1}{\sqrt{(1 + \lambda_2)}} |\beta^*|_1 \quad (7)$$

With the corrected elastic net estimates  $\hat{\beta}$  defined as

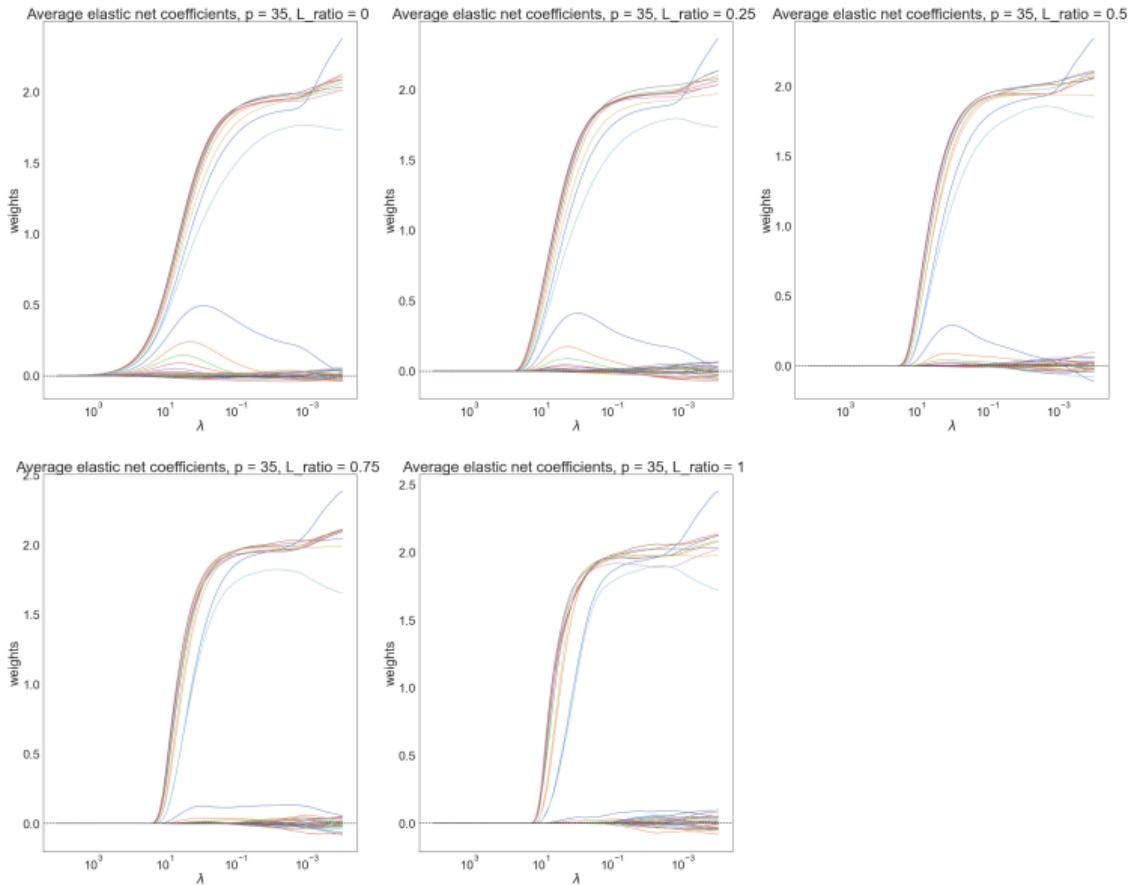
$$\hat{\beta}(\text{elastic net}) = \sqrt{(1 + \lambda_2)} \hat{\beta}^* \quad (8)$$

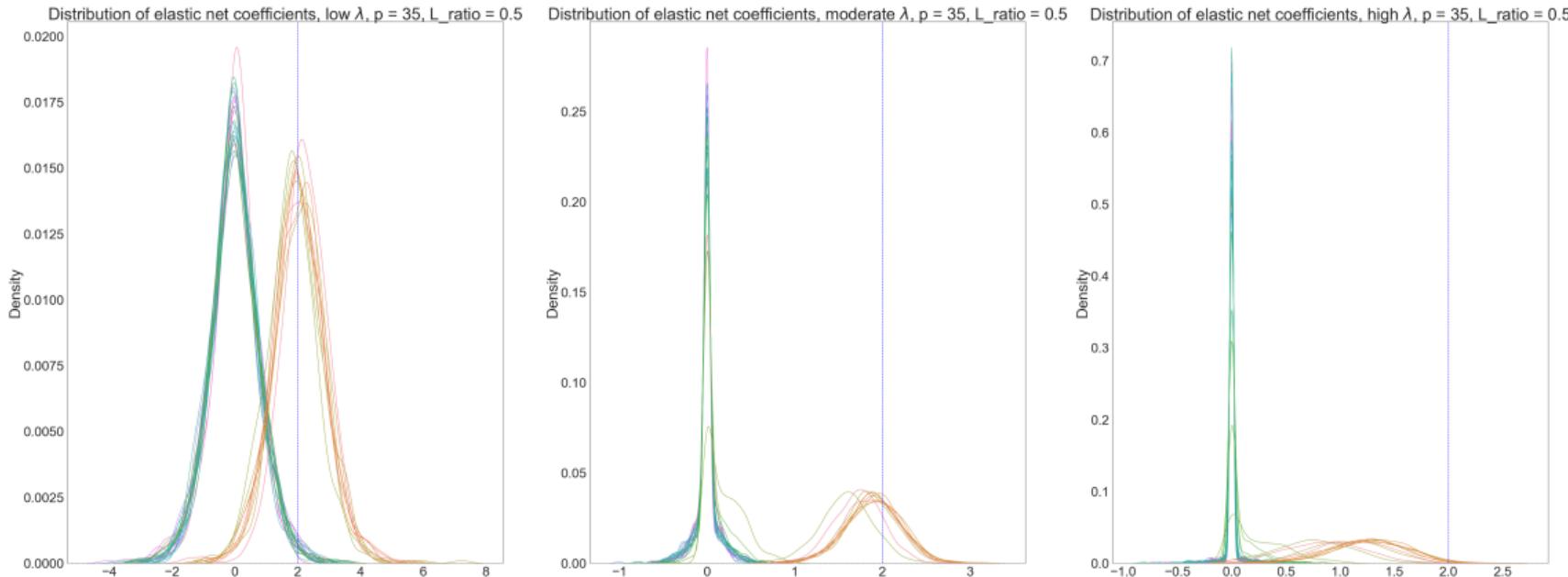
and

$$\hat{\beta}(\text{naive elastic net}) = 1 / \sqrt{(1 + \lambda_2)} \hat{\beta}^* \quad (9)$$

Finally,

$$\hat{\beta}(\text{elastic net}) = (1 + \lambda_2) \hat{\beta}(\text{naive elastic net}). \quad (10)$$





# Regularization Methods Summarized

Model	Characteristics	Drawbacks
Ridge	<ul style="list-style-type: none"> <li><math>\ell_2</math>-norm shrinkage penalty</li> <li>Performs well in the presence of multicollinearity.</li> </ul>	<ul style="list-style-type: none"> <li>Does not perform variable selection.</li> </ul>
Lasso	<ul style="list-style-type: none"> <li><math>\ell_1</math>-norm shrinkage penalty</li> <li>Performs variable selection.</li> </ul>	<ul style="list-style-type: none"> <li>Does not behave well when sparsity is very low.</li> <li>Cannot handle high correlation or grouped multicollinearity.</li> <li>Can only select up to <math>n</math> regressors when <math>p &gt; n</math></li> </ul>
Naïve Elastic Net	<ul style="list-style-type: none"> <li>convex combination of <math>\ell_1</math>-norm and <math>\ell_2</math>-norm shrinkage penalties.</li> <li>Potentially selects all <math>p</math> in the presence of grouped collinearity.</li> </ul>	<ul style="list-style-type: none"> <li>Double shrinkage does not contribute to further reduction of the variances and adds unnecessary bias.</li> <li>Works only when solution is very close to ridge or lasso.</li> </ul>
Elastic Net	<ul style="list-style-type: none"> <li>Maintains all characteristics from naïve model, but corrects for double shrinkage.</li> </ul>	

# Measuring Prediction Performance

- Selects the model that best fits the data.
- Bias-Variance Decomposition
  - We carefully follow Hastie et al. 2008.
  - Assume  $Y = f(X) + \varepsilon$ , where  $E[\varepsilon] = 0$  and  $Var(\varepsilon) = \sigma_\varepsilon^2$
  - We can derive the Mean Squared Error (MSE) out of the expected prediction error of  $\hat{f}(X)$  at a **fixed** point  $X = x_0$

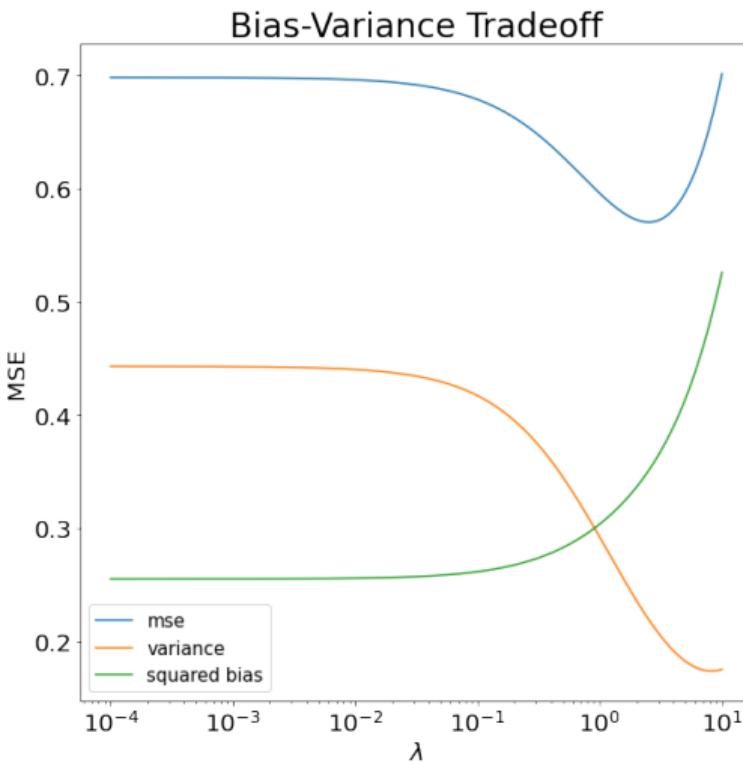
$$Err(x_0) = E[(Y - \hat{f}(x_0))^2 | X = x_0] \quad (11)$$

$$= \sigma_\varepsilon^2 + [E\hat{f}(x_0) - f(x_0)]^2 + E[\hat{f}(x_0) - E\hat{f}(x_0)]^2 \quad (12)$$

$$= \sigma_\varepsilon^2 + \text{Bias}^2(\hat{f}(x_0)) + \text{Var}(\hat{f}(x_0)) \quad (13)$$

$$= \text{Irreducible error} + \underbrace{\text{Bias}^2 + \text{Variance}}_{\text{MSE}} \quad (14)$$

# MSE as a Measurement of Model Fit



## DGP setup

$$y = X\beta + \varepsilon, \quad \text{with} \quad \varepsilon \sim \mathcal{N}_n(0, 1)$$

$X_{n \times p} \sim \mathcal{N}(0, 1)$  such that the pairwise correlation between  $x_i$  and  $x_j$  is expressed as

$$\text{corr}(i,j) = \text{correlation factor}^{|i-j|},$$

where the correlation factor takes on a value between 0 and 1.

# Minimum Test MSE (Case 1)

## Set up

- High dimensionality with varying sparsity, no multicollinearity.
- $n = 30$
- $p = 35$
- All non-zero  $\beta = 2$ .

Model	High Sparsity (10 non-zero $\beta$ 's)	Med Sparsity (20 non-zero $\beta$ 's)	Low Sparsity (35 non-zero $\beta$ 's)
	12.544	23.895	36.074
Ridge	12.544	23.895	36.074
Elnet (Naive), 0.2	10.182	21.360	40.366
Elnet (Naive), 0.5	8.760	21.687	46.961
Elnet (Naive), 0.7	7.421	22.380	52.046
Lasso	5.903	24.090	60.170

## Minimum Test MSE (Case 2)

### Set up

- Low dimensionality with varying sparsity, moderate to high pairwise correlation (corr. factor = 0.8).
- $n = 30, p = 10$
- All non-zero  $\beta = 2$ .

Model	High Sparsity	Med Sparsity	Low Sparsity
	(3 non-zero $\beta$ 's)	(7 non-zero $\beta$ 's)	(10 non-zero $\beta$ 's)
Ridge	0.570	0.270	1.140
Elnet (Naive), 0.2	0.572	0.265	1.149
Elnet (Naive), 0.5	0.577	0.246	1.164
Elnet (Naive), 0.7	0.552	0.218	1.166
Lasso	0.607	0.298	1.166

## Minimum Test MSE (Case 3)

### Set up

- Low dimensionality with high sparsity, varying degrees of pairwise correlation.
- $n = 20$
- $p = 8$
- $\beta = (3, 1.5, 0, 0, 2, 0, 0, 0)$

Model	Low Pairwise Corr.	Med Pairwise Corr.	High Pairwise Corr.
	(corr. factor = 0.1)	(corr. factor = 0.3)	(corr. factor = 0.7)
Ridge	1.238	2.071	0.905
Elnet (Naive), 0.2	1.160	2.034	0.889
Elnet (Naive), 0.5	0.979	1.897	0.864
Elnet (Naive), 0.7	0.828	1.724	0.851
Lasso	0.600	1.340	0.829

## Minimum Test MSE (Case 4)

### Set up

- Low dimensionality with no sparsity, varying degrees of pairwise correlation.
- $n = 20$
- $p = 8$
- All  $\beta = 0.85$

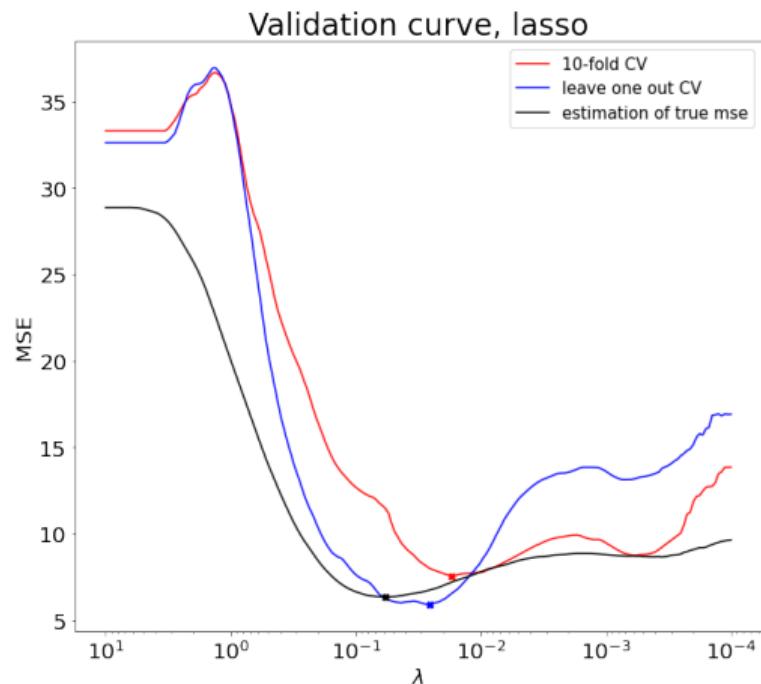
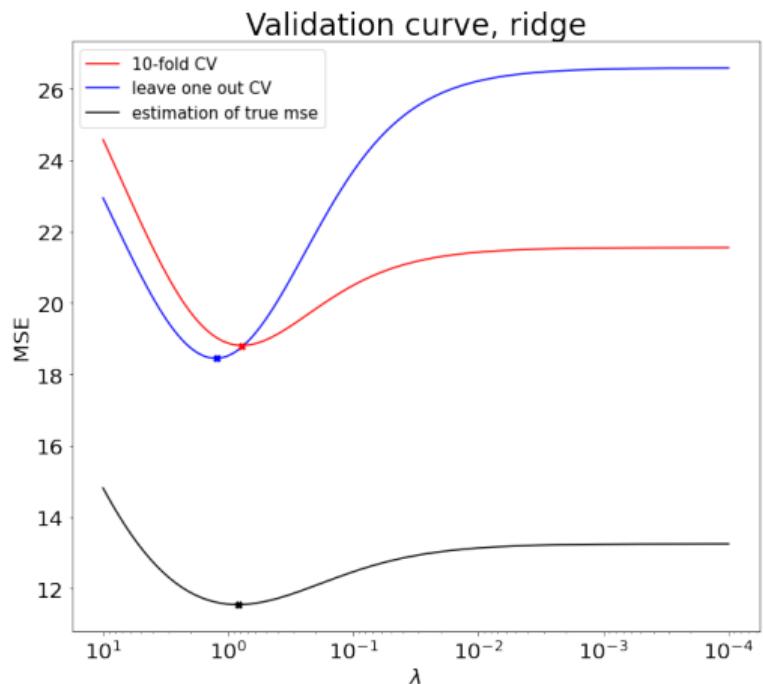
Model	Low Pairwise Corr. (corr. factor = 0.1)	Med Pairwise Corr. (corr. factor = 0.3)	High Pairwise Corr. (corr. factor = 0.7)
Ridge	2.471	0.521	0.736
Elnet (Naive), 0.2	2.506	0.585	0.772
Elnet (Naive), 0.5	2.542	0.882	0.831
Elnet (Naive), 0.7	2.548	1.124	0.867
Lasso	2.548	1.532	0.906

# Minimum Test MSE (Case 5)

## Set up

- High dimensionality with high sparsity, varying degrees of pairwise correlation.
- $n = 30$
- $p = 35$
- 10 non-zero  $\beta = 2$

Model	Low Pairwise Corr. (corr. factor = 0.1)	Med Pairwise Corr. (corr. factor = 0.3)	High Pairwise Corr. (corr. factor = 0.7)
Ridge	9.412	15.714	3.263
Elnet (Naive), 0.2	7.932	11.092	2.659
Elnet (Naive), 0.5	6.049	10.505	1.882
Elnet (Naive), 0.7	5.122	8.570	1.589
Lasso	4.079	6.575	1.716

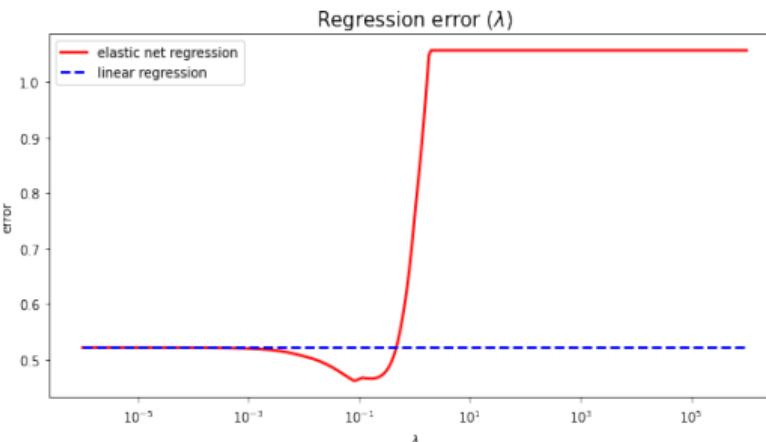
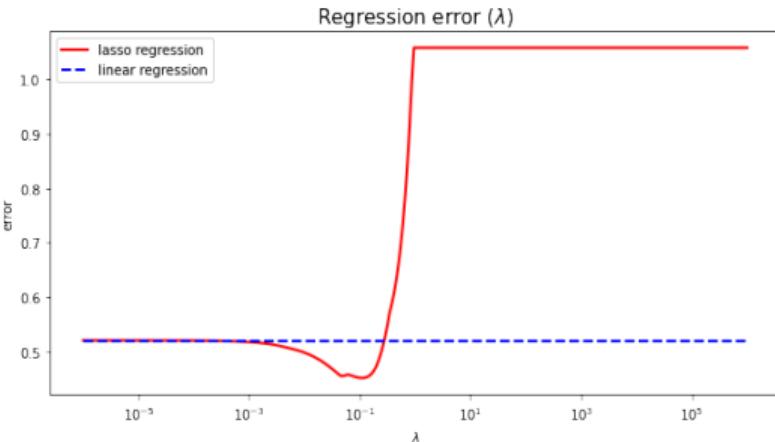
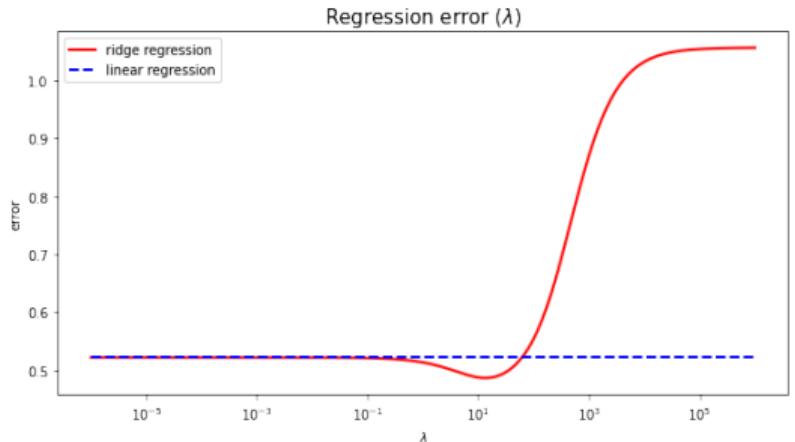


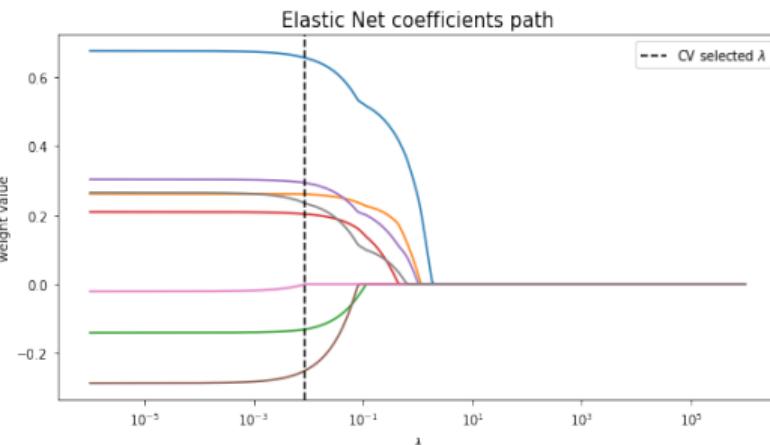
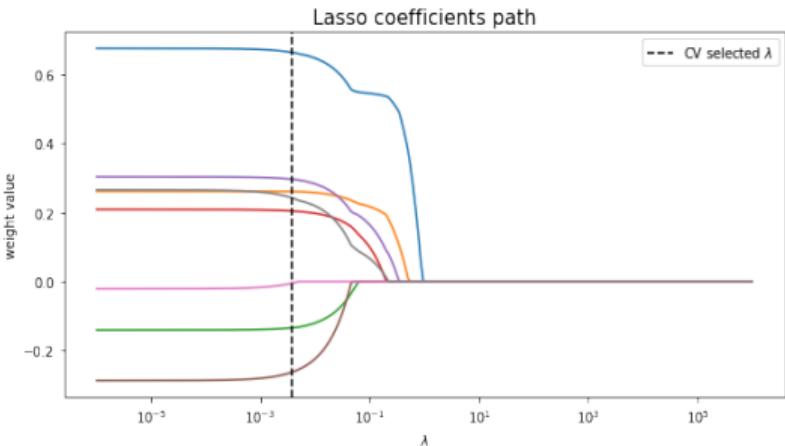
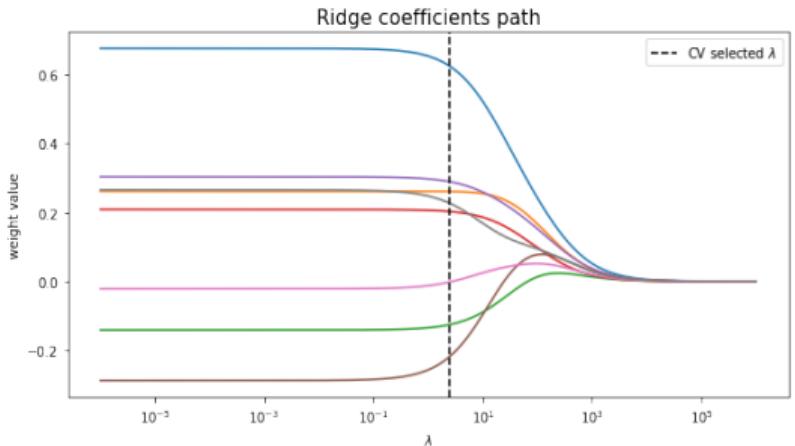
# Prostate Cancer Data Set

- What is the relationship between log prostate specific antigen (lpsa) levels, which is elevated in men with prostate cancer, and other clinical measures ( $p = 8$ )
- $n = 97$
- For comparability we use the same train-test split as Zou et al. 2005
- High correlation between
  - lcp (*log of capsular penetration*) and lcavol (*log cancer volume*)
  - svi (*seminal vesicle invasion*) and lcp
  - pgg45 (*percent of Gleason score 4 or 5*) and gleason (*Gleason score*)
  - pgg45 and lcp

# Correlation Matrix

	Ipsa	Icavol	Iweight	age	lbph	svi	Icp	gleason	pgg45
Ipsa	1.000000	0.734460	0.433319	0.169593	0.179809	0.566218	0.548813	0.368987	0.422316
Icavol	0.734460	1.000000	0.280521	0.225000	0.027350	0.538845	0.675310	0.432417	0.433652
Iweight	0.433319	0.280521	1.000000	0.347969	0.442264	0.155385	0.164537	0.056882	0.107354
age	0.169593	0.225000	0.347969	1.000000	0.350186	0.117658	0.127668	0.268892	0.276112
lbph	0.179809	0.027350	0.442264	0.350186	1.000000	-0.085843	-0.006999	0.077820	0.078460
svi	0.566218	0.538845	0.155385	0.117658	-0.085843	1.000000	0.673111	0.320412	0.457648
Icp	0.548813	0.675310	0.164537	0.127668	-0.006999	0.673111	1.000000	0.514830	0.631528
gleason	0.368987	0.432417	0.056882	0.268892	0.077820	0.320412	0.514830	1.000000	0.751905
pgg45	0.422316	0.433652	0.107354	0.276112	0.078460	0.457648	0.631528	0.751905	1.000000





# Model Selection

Model	Test MSE (with 10-fold CV)	Test MSE (AIC)	Variable Selection
Linear	0.5213	-	All
Ridge	0.5043	-	All
Lasso	0.5112	-	All
Elastic Net (Naive)	0.5043	-	All
LARS	0.5084	0.5033	lcavol, lweight, lbph, svi, pgg45

## Takeaway

- Naive elastic net performs as well as ridge, also reported by Zou et al. 2005
- Lasso performs best if LARS algorithm is implemented and AIC is used for tuning  $\lambda$ . It performs variable selection as in Zou et al. 2005.

## Summary: best model based on simulations

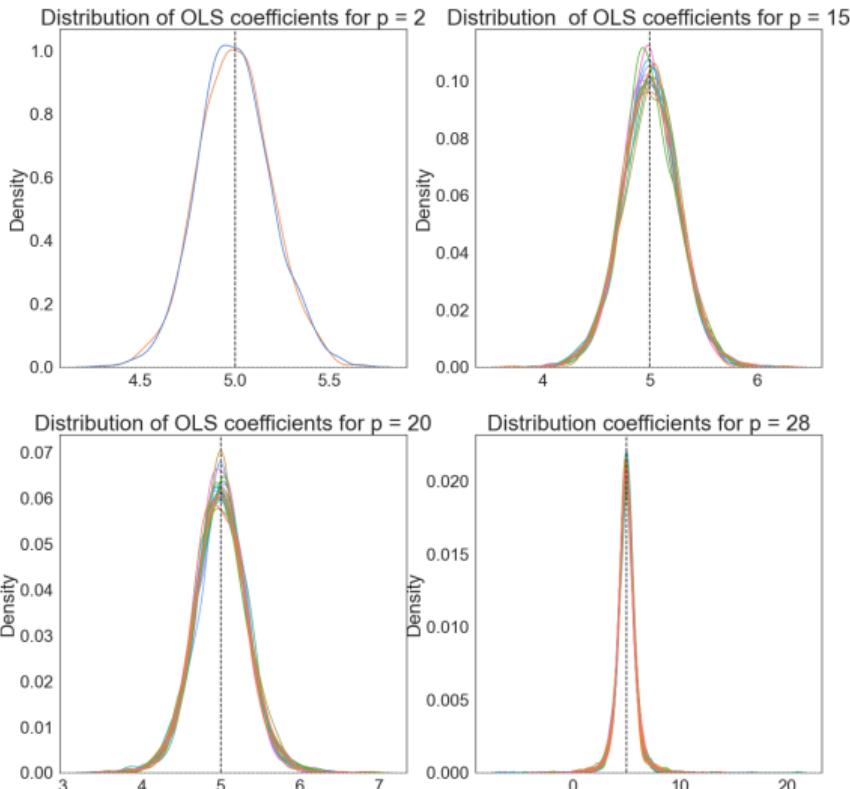
		Regressors' correlation	
		High	Low
Coefficients' sparsity	High	Elastic Net (high-dimensionality) Lasso (low-dimensionality)	Lasso
	Low	Ridge	Ridge (high-dimensionality) OLS (low-dimensionality)

- The data application is in line with *Case 2* of our simulation study.

**Thank you for your attention!**



# Appendix



# Appendix

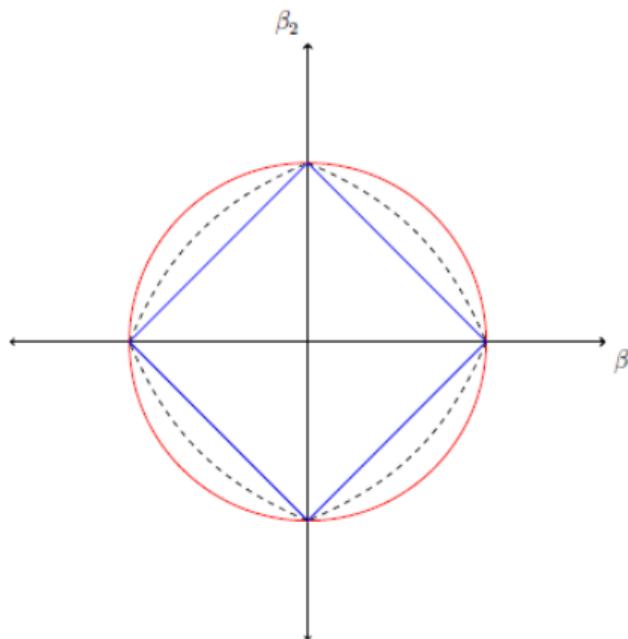
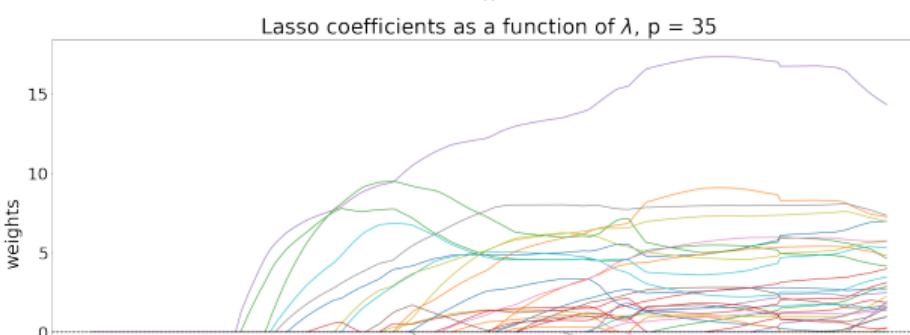
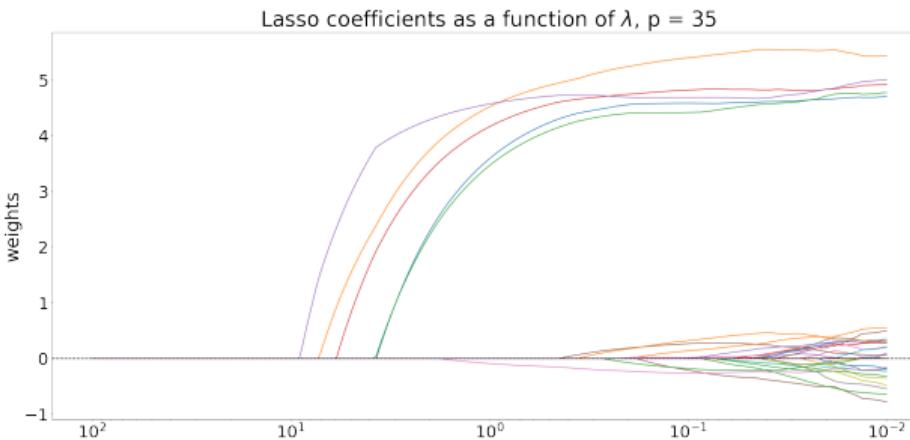


Figure 2: constraint region for Ridge (red) L2-norm, Lasso (blue) L1-norm, Elastic Net (dashed black)

# Appendix



## Appendix

- We used a *10-fold cross validation* approach to find the optimal value of the tuning parameter.

$$CV_{(k)} = \frac{1}{k} \sum_{i=1}^k MSE_i \quad (15)$$

- We identify the model with the lowest test error.
- We identify the  $\lambda$  value that minimizes the MSE and then plug it back into the model of interest.

# Appendix

## Least-angle regression (LARS) (Efron et al. 2003)

- Provides an extremely efficient algorithm for computing the entire lasso path.
- Lasso is a variant of the LARS procedure.
- LARS algorithm:
  - 1 Standardize predictors.
  - 2 Find predictor  $x_j$  that is most correlated with  $y$ .
  - 3 Move  $\beta_j$  from 0 towards least-squares coefficient until  $x_k$  has as much correlation with the current residual as does  $x_j$ .
  - 4 Move  $\beta_j$  and  $\beta_k$  in the direction of the joint least squares coefficient of the current residual on  $(x_j, x_k)$ , until  $x_l$  has as much correlation with the current residual.
  - 5 Continue until all  $p$  have been entered.
  - 6 *If a non-zero coefficient hits zero, drop the corresponding variable from the active set of variables and recompute the current joint least squares direction.*

# References |

- Aheto, J. M. K. et al. (2021). "A predictive model, and predictors of under-five child malaria prevalence in Ghana: How do LASSO, Ridge and Elastic net regression approaches compare?" In: *Preventive Medicine Reports* 23, p. 101475.
- Efron, B. et al. (2003). "Least Angle Regression". In: *The Annals of Statistics* 32.2, pp. 407–451.
- Hastie, T., R. Tibshirani, and J. Friedman (2008). *Elements of Statistical Learning*. Springer Series in Statistics.
- Hastie, T., R. Tibshirani, and M. Wainwright (2019). *Statistical learning with sparsity: the lasso and generalizations*. Chapman and Hall/CRC.
- Hoerl, A. E. and R. W. Kennard (1970a). "Ridge regression: applications to nonorthogonal problems". In: *Technometrics* 12.1, pp. 69–82.
- (1970b). "Ridge regression: Biased estimation for nonorthogonal problems". In: *Technometrics* 12.1, pp. 55–67.
- James, G. et al. (2013). *An introduction to statistical learning*. Vol. 112. Springer.
- Kidwell, J. S. and L. H. Brown (1982). "Ridge regression as a technique for analyzing models with multicollinearity". In: *Journal of Marriage and the Family*, pp. 287–299.
- Marquardt, D. W. and R. D. Snee (1975). "Ridge regression in practice". In: *The American Statistician* 29.1, pp. 3–20.
- Melkumova, L. and S. Y. Shatskikh (2017). "Comparing Ridge and LASSO estimators for data analysis". In: *Procedia engineering* 201, pp. 746–755.
- Snee, R. D. (1973). "Some aspects of nonorthogonal data analysis: Part I. Developing prediction equations". In: *Journal of Quality Technology* 5.2, pp. 67–79.
- Tibshirani, R. (1996). "Regression shrinkage and selection via the lasso". In: *Journal of the Royal Statistical Society: Series B (Methodological)* 58.1, pp. 267–288.

## References II

Zou, H. and T. Hastie (2005). "Regularization and variable selection via the elastic net". In: *Journal of the royal statistical society: series B (statistical methodology)* 67.2, pp. 301–320.