

Regularized regression

Readings for today

- Chapter 6: Linear model selection and regularization. James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). An introduction to statistical learning: with applications in R (Vol. 6). New York: Springer

Topics

1. Ridge regression

2. LASSO

3. Elastic net

Ridge regression

Dimensionality

Dimensionality of a model: $n \times p$

As $n \rightarrow p$, dimensionality increases & model variance increases

$$\begin{pmatrix} x_{1,1} \\ \dots \\ x_{n,1} \end{pmatrix} \rightarrow \begin{pmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,15} \\ \dots & & & \\ x_{n,1} & x_{n,2} & \dots & x_{n,15} \end{pmatrix} \rightarrow \uparrow \text{model flexibility}$$

How do you select (or remove) variables in the fitting process itself?

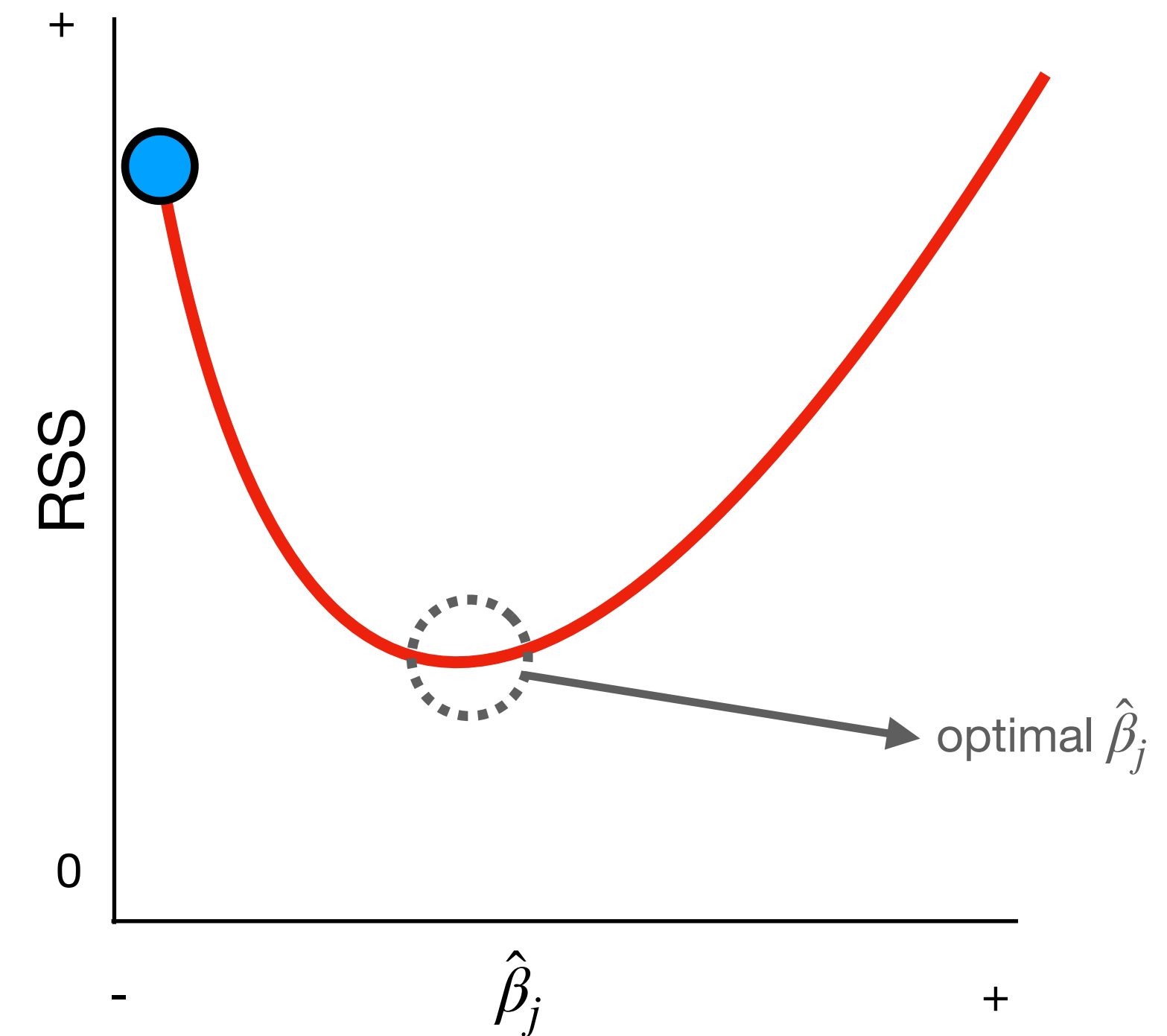
OLS Objective

Residual Sums of Squares (RSS):

$$\begin{aligned} RSS &= e_1^2 + \dots + e_n^2 \\ &= (y_1 - \sum_{j=1}^p \hat{\beta}_j x_{1,j})^2, \dots, \sum_{j=1}^p \hat{\beta}_j x_{n,j})^2 \\ &= \sum_{i=1}^n (y_i - \sum_{j=1}^p \hat{\beta}_j x_{i,j})^2 \end{aligned}$$

Objective: $\min(\log(L)) = \min(\sum_{i=1}^n (y_i - \hat{y}_i)^2)$

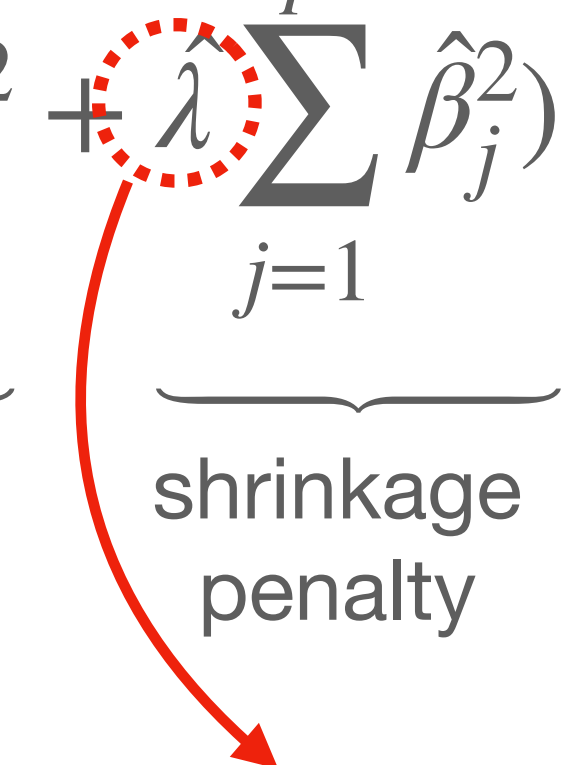
$$= \min(\sum_{i=1}^n (y_i - \sum_{j=1}^p \hat{\beta}_j x_{i,j})^2)$$



Ridge regression

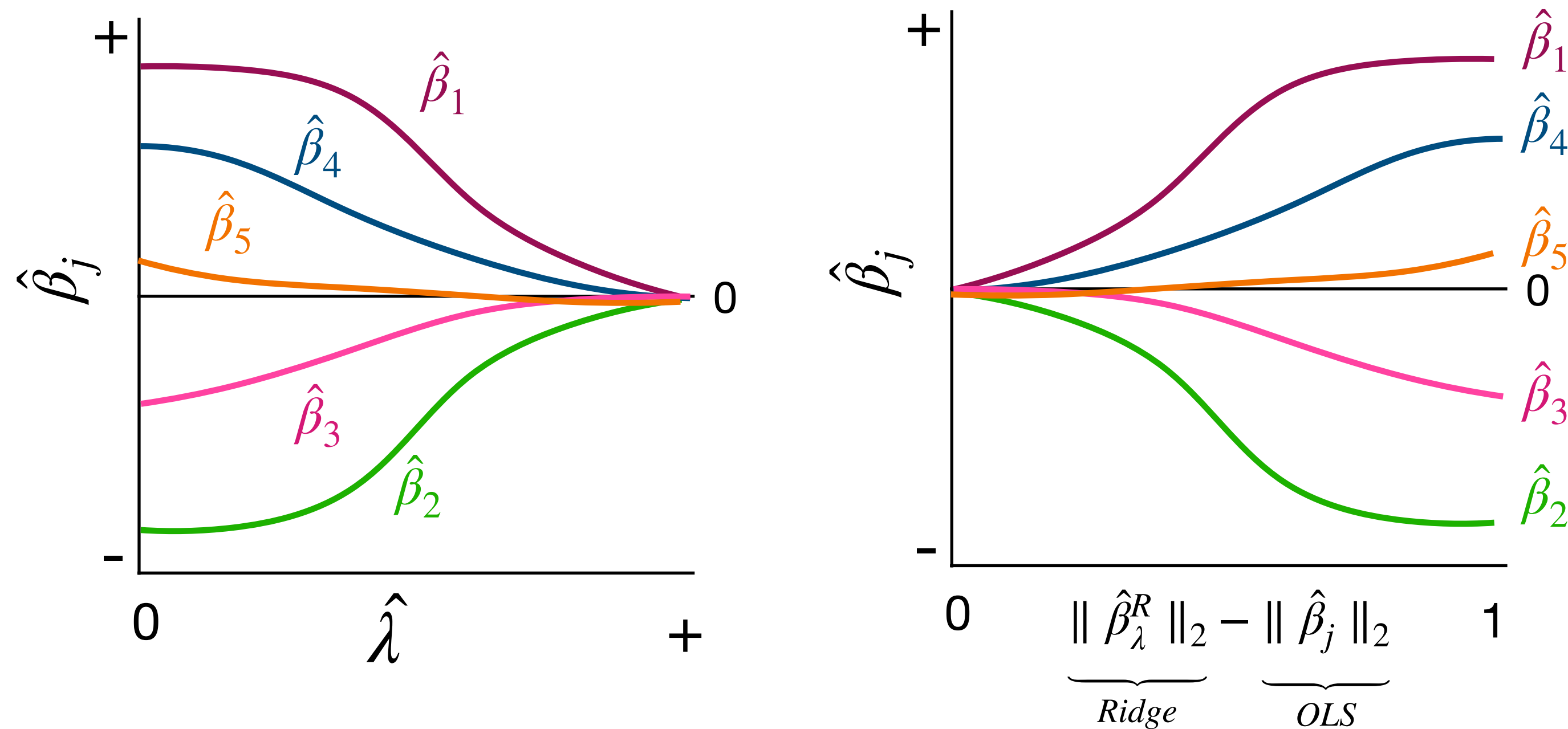
Goal: Penalize “weak” effects so they have minimal impact on your model.

Objective: $\min(\log(L)) = \min(\underbrace{\sum_{i=1}^n (y_i - \sum_{j=1}^p \hat{\beta}_j x_{i,j})^2}_{\text{minimized residual squared error}} + \underbrace{\hat{\lambda} \sum_{j=1}^p \hat{\beta}_j^2}_{\text{shrinkage penalty}})$



$\hat{\lambda}$ = sparsity parameter

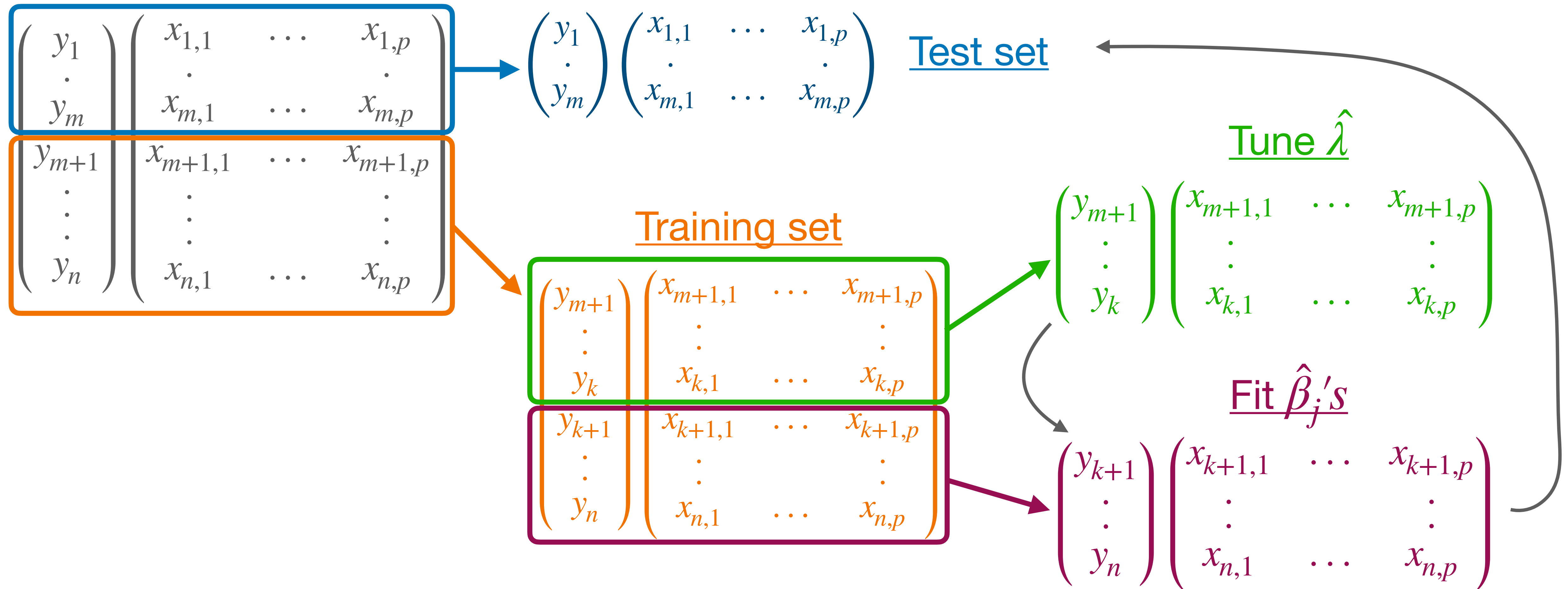
Effects on coefficients



- Ridge will *always* return all p predictors.
- $\uparrow \lambda$ results in weaker $\hat{\beta}_j$'s.
- Ridge manages the bias-variance tradeoff by reducing the influence of weak predictor variables.

Does *not* select features!

Finding λ



$\hat{\lambda}$ is a free parameter that needs to be fit

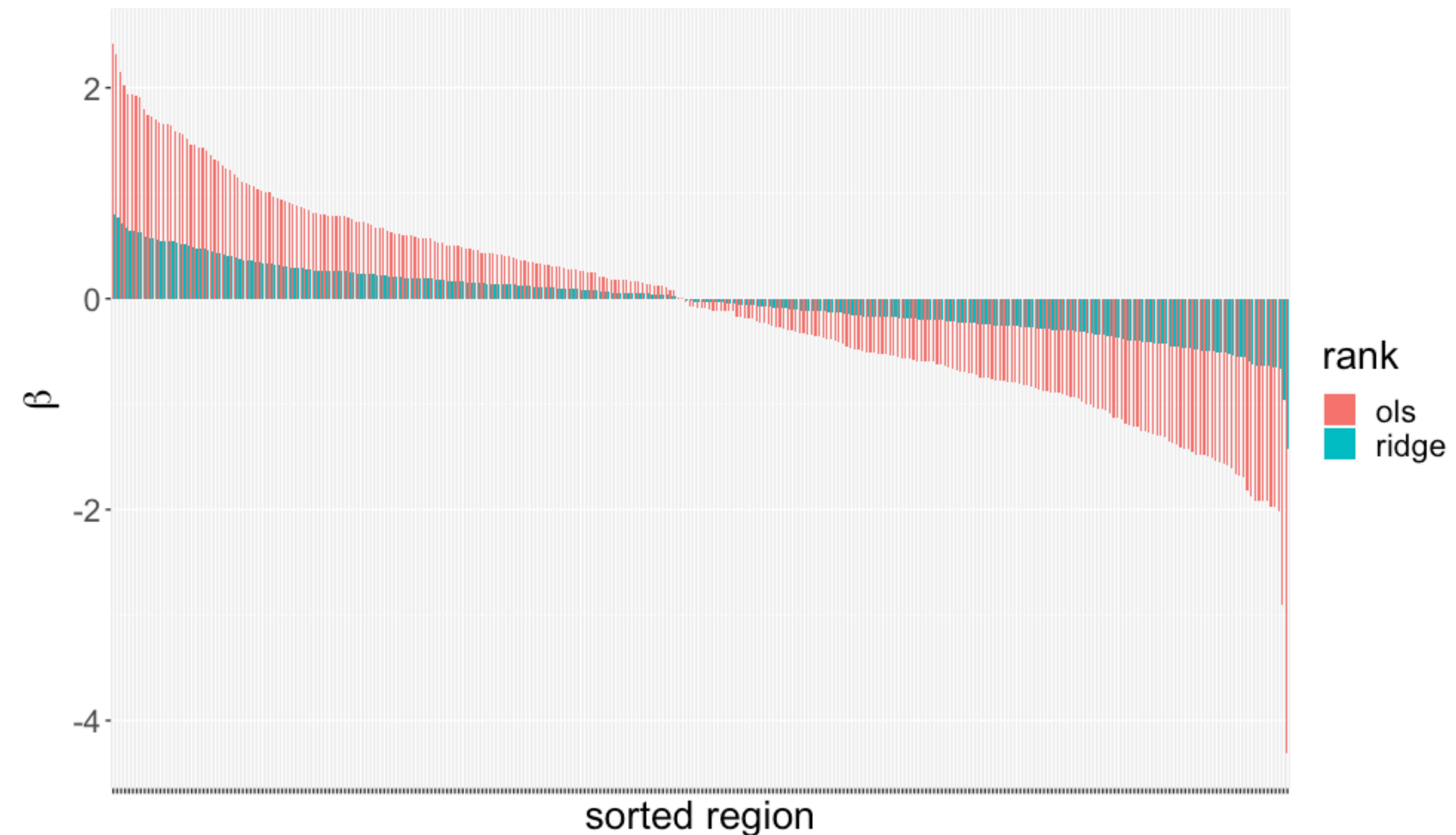
Example: high variance model

Q: What brain areas predict reaction times across trials?

$$Y_{rt,t} = \hat{\beta}_0 + \sum_{j=1}^R \hat{\beta}_j X_{t,j} \quad \begin{array}{l} R = 300 \\ n = 500 \end{array}$$

$Y_{rt,t}$: reaction time on trial t

$X_{j,t}$: fMRI response of region j on trial t



The pros and cons of Ridge regression

Advantages:

- Manages the bias-variance tradeoff well when $\frac{p}{n}$ is very high.
- Best solution for high variance models.

Disadvantages:


- Does *not* actually do feature selection (simply pushes $\hat{\beta}_j$'s close to zero).
- Is not scale invariant
 - OLS: $cX \rightarrow c\hat{\beta}$
 - Ridge: $cX \neq c\hat{\beta}$

LASSO

Least Absolute Shrinkage & Selection Operator (LASSO)

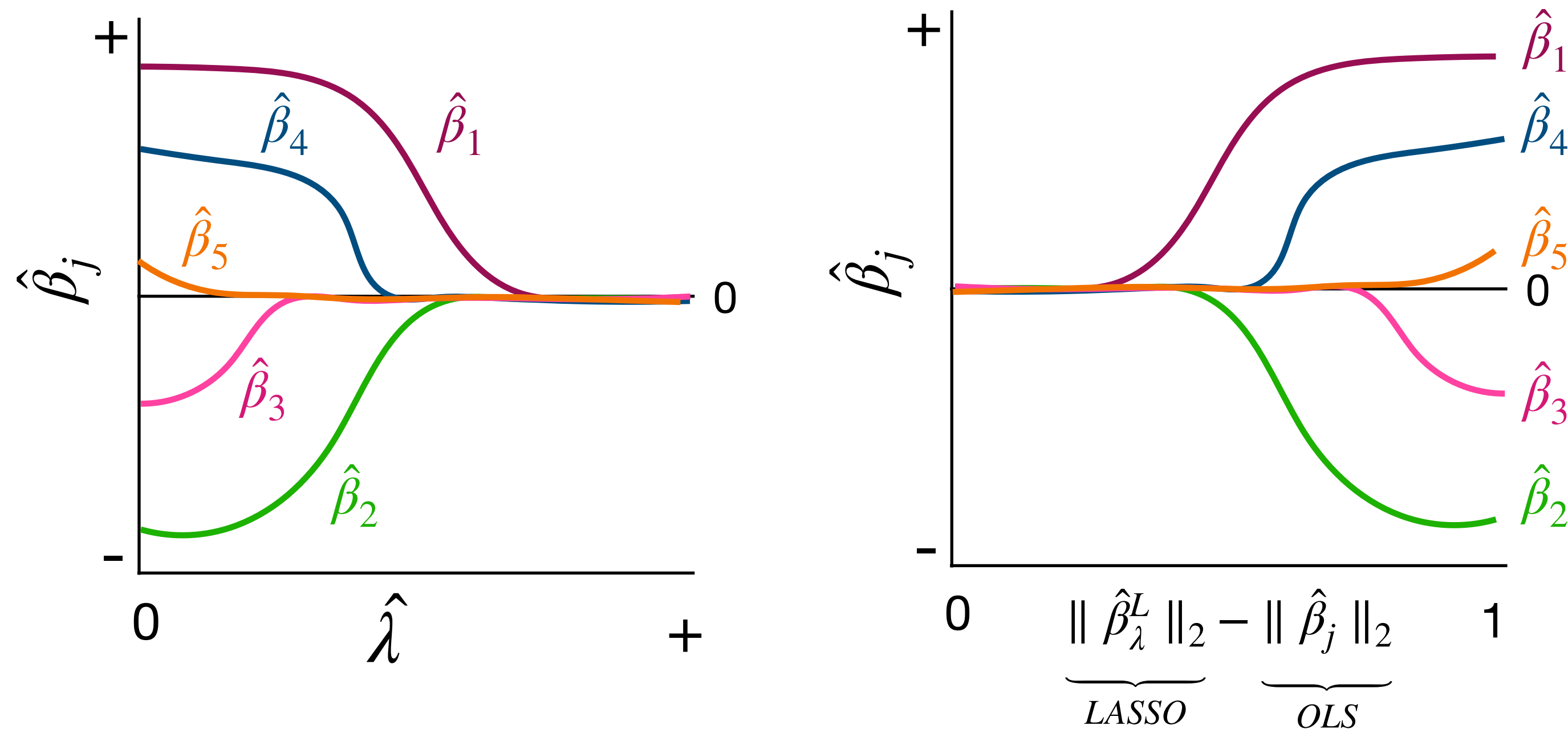
Goal: Remove “weak” effects altogether

Objective: $\min(\log(L)) = \min(\underbrace{\sum_{i=1}^n (y_i - \sum_{j=1}^p \hat{\beta}_j x_{i,j})^2}_{\text{minimized residual squared error}} + \underbrace{\hat{\lambda} \sum_{j=1}^p |\hat{\beta}_j|}_{\text{shrinkage penalty}})$



$\hat{\beta}_j$'s can now shrink to 0.

Effects on coefficients




- LASSO returns $< p$ predictors with $\lambda > 0$.
- $\uparrow \lambda$ results in more $\hat{\beta}_j$'s = 0.
- LASSO does feature selection in the model fit process by setting coefficients to zero.

Symmetry: Ridge vs. LASSO

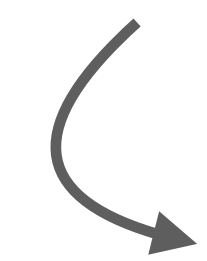
Assume:

$$X = I = \begin{pmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & & & & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{pmatrix} \quad p = n$$

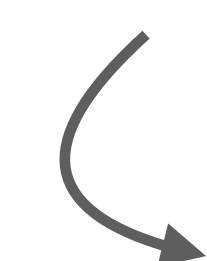
OLS:

$$\min(\log(L)) = \min\left(\sum_{i=1}^n (y_i - \sum_{j=1}^p \hat{\beta}_j x_{i,j})^2\right)$$

$$\hat{\beta}_i = y_i$$

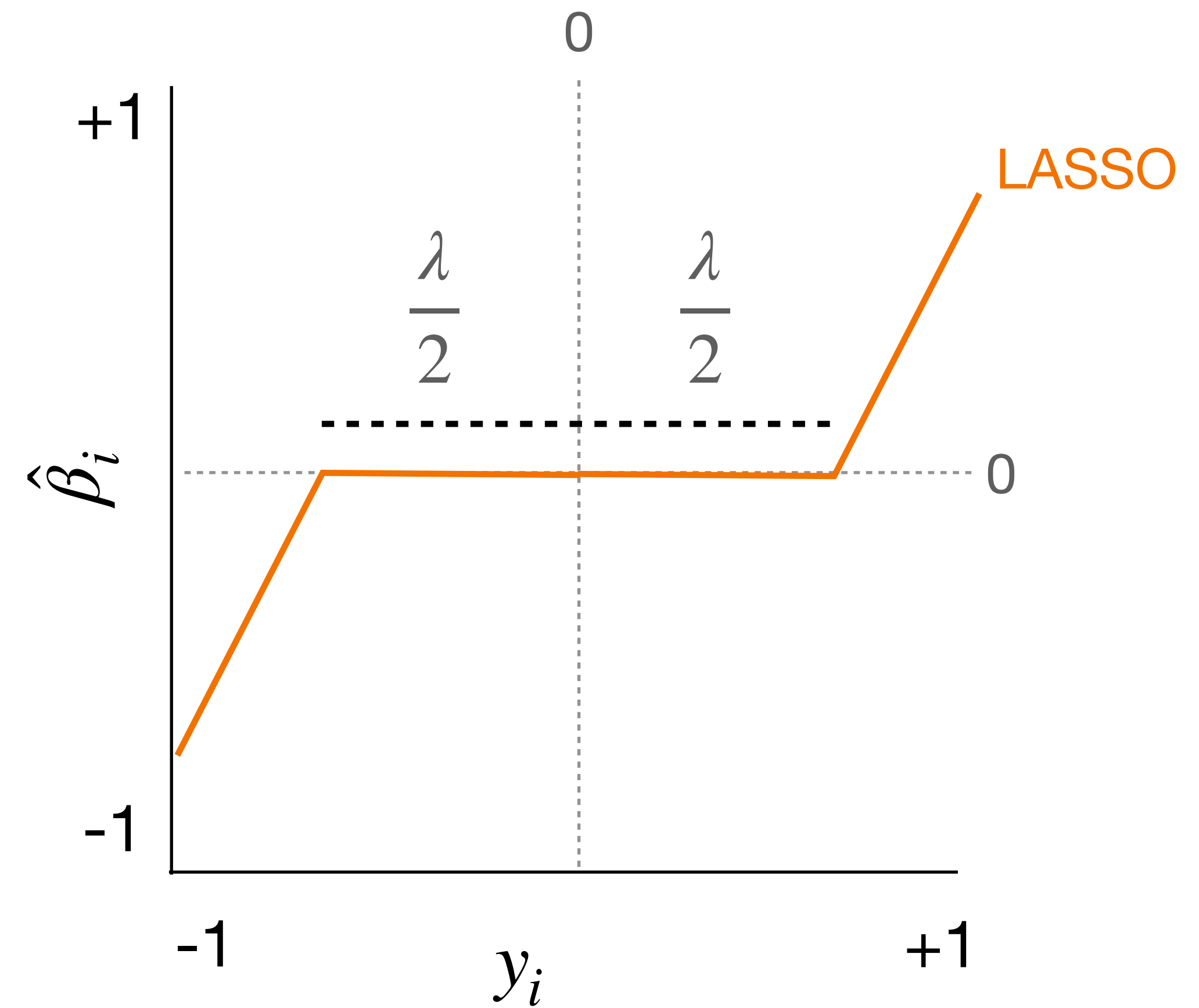
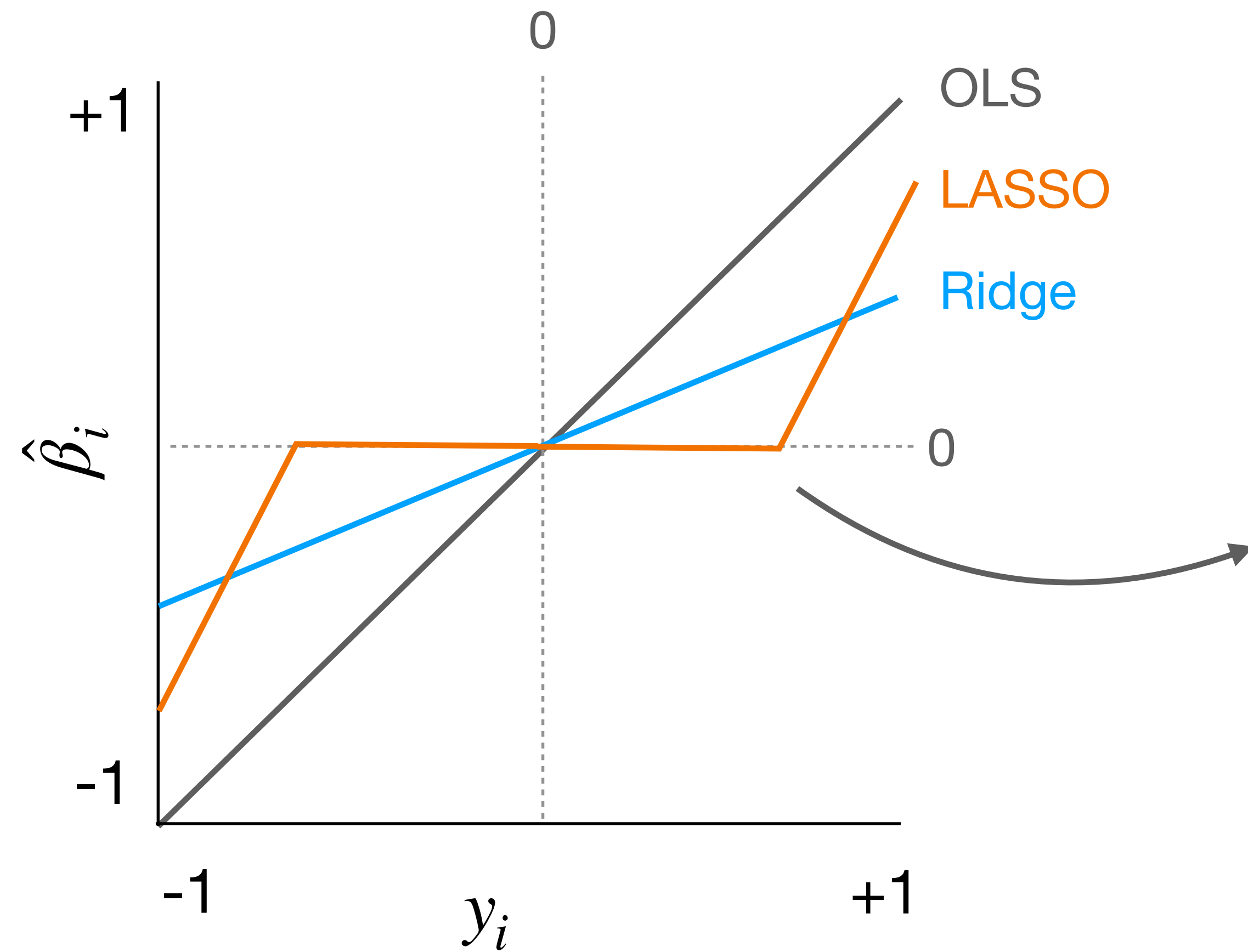
Ridge:

$$\min(\log(L)) = \min\left(\sum_{i=1}^n (y_i - \sum_{j=1}^p \hat{\beta}_j x_{i,j})^2 + \hat{\lambda} \sum_{j=1}^p \hat{\beta}_j^2\right)$$

$$\hat{\beta}_i^R = \frac{y_i}{1 + \hat{\lambda}}$$

LASSO:

$$\min(\log(L)) = \min\left(\sum_{i=1}^n (y_i - \sum_{j=1}^p \hat{\beta}_j x_{i,j})^2 + \hat{\lambda} \sum_{j=1}^p |\hat{\beta}_j|\right)$$

$$\hat{\beta}_i^L = \begin{cases} y_i - \frac{\hat{\lambda}}{2}, & \text{if } y_i > \frac{\hat{\lambda}}{2} \\ 0, & \text{if } |y_i| \leq \frac{\hat{\lambda}}{2} \\ y_i + \frac{\hat{\lambda}}{2}, & \text{if } y_i < -\frac{\hat{\lambda}}{2} \end{cases}$$

Symmetry: Ridge vs. LASSO



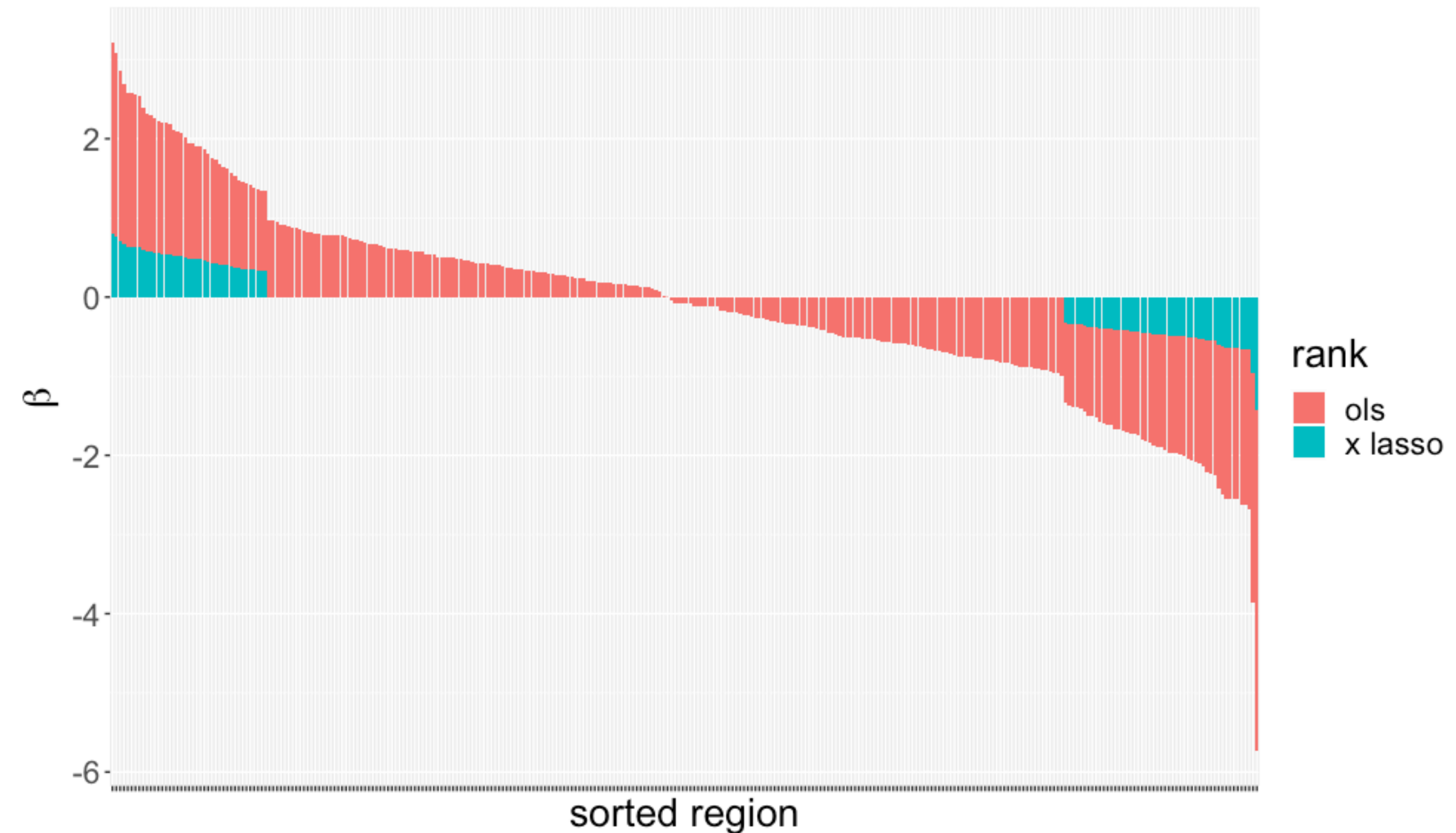
Example: high variance model

Q: What brain areas predict reaction times across trials?

$$Y_{rt,t} = \hat{\beta}_0 + \sum_{j=1}^R \hat{\beta}_j X_{t,j} \quad \begin{array}{l} R = 300 \\ n = 500 \end{array}$$

$Y_{rt,t}$: reaction time on trial t

$X_{j,t}$: fMRI response of region j on trial t

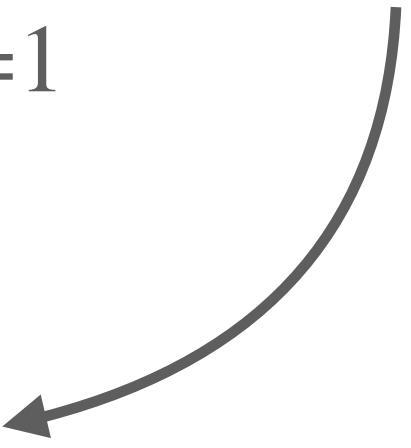


Best subset as LASSO

LASSO can implement the best subset algorithm with a change in the sparsity penalty.

Best subset penalty:

$$\min(\log(L)) = \min\left(\sum_{i=1}^n (y_i - \sum_{j=1}^p \hat{\beta}_j x_{i,j})^2 + \hat{\lambda} \sum_{j=1}^p I(\hat{\beta}_j \neq 0)\right)$$

$$I(\hat{\beta}_j \neq 0) = \begin{cases} 1, & \text{if } \hat{\beta}_j \neq 0 \\ 0, & \text{otherwise} \end{cases}$$


Ridge vs. LASSO vs. Subset

1. LASSO picks the “best” of any correlated/clustered set of predictor variables, X .
2. Ridge picks the strongest subset of correlated/clustered predictor variables, X .
3. LASSO is more conservative than Ridge or best subset selection.

Elastic net

Elastic net

Goal: Little bit of ridge and little bit of LASSO

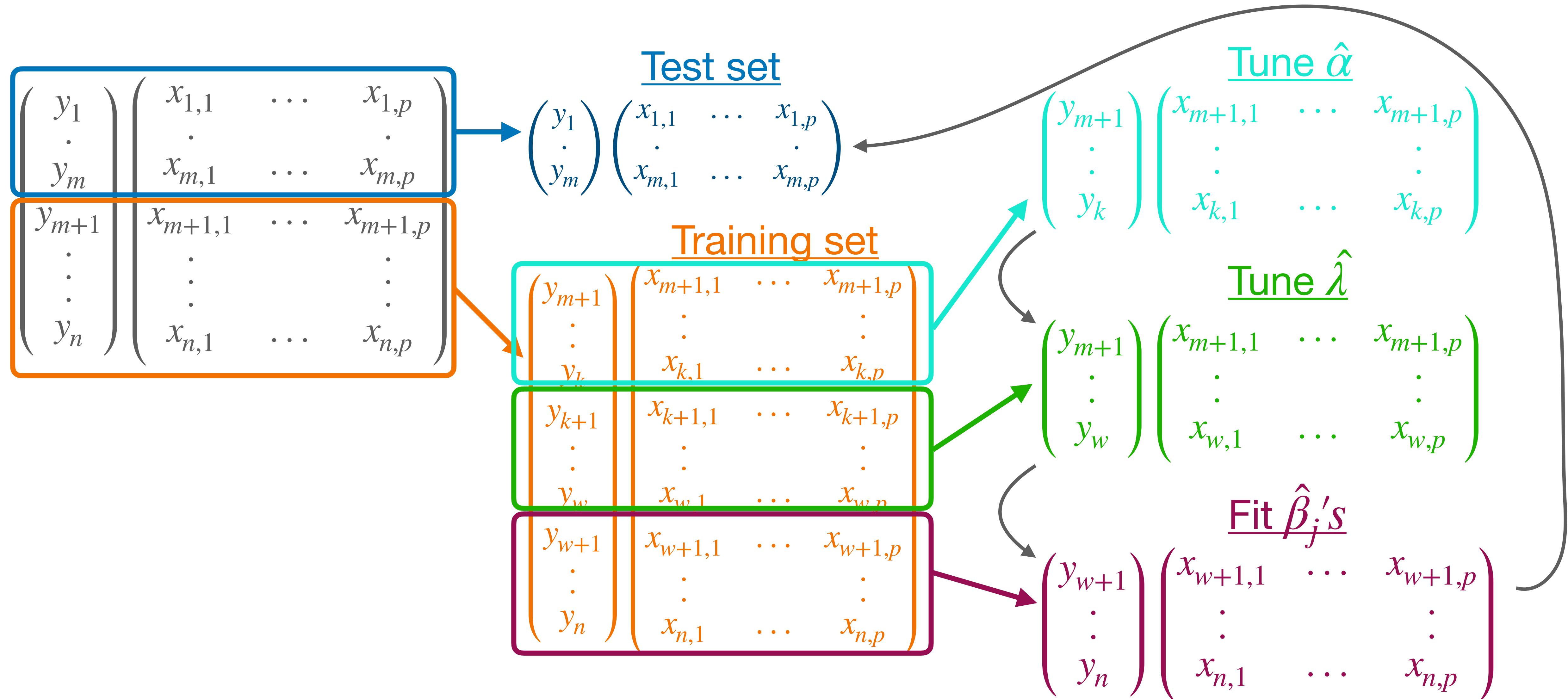
Objective: $\min(\log(L)) = \min(\underbrace{\sum_{i=1}^n (y_i - \sum_{j=1}^p \hat{\beta}_j x_{i,j})^2}_{\text{minimized residual squared error}} + \underbrace{\hat{\alpha}\hat{\lambda} \sum_{j=1}^p \hat{\beta}_j^2}_{\text{partial ridge penalty}} + \underbrace{(1 - \hat{\alpha})\hat{\lambda} \sum_{j=1}^p |\hat{\beta}_j|}_{\text{partial LASSO penalty}})$

Ridge: $\alpha = 1$

$\alpha = \text{mixing parameter} \rightarrow$ LASSO: $\alpha = 0$

Elastic Net: $0 < \alpha < 1$

Finding α



Take home message

- Regularized regression models are a fast, computationally efficient way of managing model complexity by implementing feature selection (or sparsity) in the model fitting process itself.