



S&DS 365 / 565
Data Mining and Machine Learning

Regularization

Yale

Regularization helps control (regulate) parameters.

help robustness, prediction accuracy, often L2 norm

General form

$$\hat{\beta} \in \arg \min_{\beta} \underbrace{\mathcal{L}(\beta, \{Z_i\}_{i=1}^n)}_{\text{Loss Function}} + \lambda \underbrace{R(\beta)}_{\text{Regularizer}}$$

↓
param you identified
any data {
 (x_i, y_i)
 (x_i)
} no label supervised learning
 unsupervised learning

- Loss Function measures fit to data $\{Z_i\}$
the param β you identified
complexity of
- Regularizer controls parameters, fits our prior beliefs of the data
Later talk about Bayes, priors fit this general form too
- Trade-off Loss and Regularization via regularization parameter λ
 λ control how much you emphasize on regularizer
 λ is a hyper param for your problem, often use cross validation to pick optimal λ
$$\begin{cases} \lambda = 0 & \hat{\beta} \in \arg \min \text{Loss} \\ \lambda \uparrow & \hat{\beta} \in \arg \min \text{more } R(\beta) \end{cases}$$

Choices for **regularizer**: ℓ_p norms convex for $p \geq 1$

$$\beta \in \mathbb{R}^d = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_d \end{bmatrix}$$

Focus on ℓ_2 & ℓ_1 in

$$\|\beta\|_p = \left(\sum_{j=1}^d |\beta_{(j)}|^p \right)^{\frac{1}{p}}$$

discussion

ℓ_2 norm squared

$$\frac{1}{2} \|\beta\|_2^2$$

ℓ_1 norm

$$\|\beta\|_1 = \sum_{j=1}^d |\beta_{(j)}|$$

Choices for **losses** : least squares, logistic, any GLM, etc...
any loss function

① Specific example: Ridge Regression

$$\hat{\beta} \in \arg \min_{\beta} \frac{1}{2n} \|X\beta - y\|_2^2 + \frac{\lambda}{2} \|\beta\|_2^2$$

least square regression regularizer ℓ_2 -norm

Result is that the coefficients get shrunken towards 0 and standard errors of coefficients are much lower. Regularization often called **shrinkage** → n param are shrinking as $\lambda \uparrow$

② Specific example: Lasso: combines shrinkage and variable selection. (Least absolute shrinkage and selection operator) (feature / attribute selection)

$$\hat{\beta} \in \arg \min_{\beta} \frac{1}{2n} \|X\beta - y\|_2^2 + \lambda \|\beta\|_1$$

\downarrow
 ℓ_1 -norm

Technically this is Lasso

original version of Lasso - more constrained

ℓ_1 -norm is selecting variables, it set some coefficient β to be 0
good for interpretability: this 5 params are truly needed for

$$\hat{\beta} \in \arg \min_{\beta} \frac{1}{2n} \|X\beta - y\|_2^2 \text{ such that } \|\beta\|_1 \leq R$$

prediction

because convexity

→ one solver called Lasso will give solution path for all choices of R it will give corresponding $\hat{\beta}$, if $\# \text{ data} \ll \# \text{ features}$

avoid overfitting

But they are equivalent with respective choices of λ and R .

use cross validation to pick optimal λ and R

eg. 500 patients CMS vs. 30,000 genes

Ridge regression

Ridge regression has a closed form solution

$$\hat{\beta}^{(\lambda)} \in \arg \min_{\beta} \frac{1}{2n} \|X\beta - y\|_2^2 + \frac{\lambda}{2} \|\beta\|_2^2$$

Closed-form solution: HW! ⁵

similar to Linear regression: take gradient, set to 0

General loss: gradient descent (neural net people call it weight decay)

$$\hat{\beta} \in \arg \min_{\beta} \underbrace{\mathcal{L}(\beta, \{Z_i\}_{i=1}^n)}_{\text{Loss Function}} + \lambda \underbrace{R(\beta)}_{\text{Regularizer}}$$

Take $R(\beta) = \frac{1}{2} \|\beta\|_2^2$ l₂-norm (penalty)

$$= \beta_k - \eta_k \nabla_{\beta} \mathcal{L}(\beta_k, \{Z_i\}_{i=1}^n) = \beta_k - \eta_k (\nabla_{\beta} \mathcal{L}(\beta_k) + \lambda \beta_k)$$

$$\beta_{k+1} = (1 - \eta_k \lambda) \beta_k - \eta_k \nabla \mathcal{L}(\beta_k, \{Z_i\}_{i=1}^n)$$

λ for ridge regression

- Each λ results in a different set of coefficients $\hat{\beta}^{(\lambda)}$.
- λ controls the amount of shrinkage (bias-variance tradeoff).

▶ λ near 0 implies $\hat{\beta}^{(\lambda)}$ is near least squares estimate $\hat{\beta}$.

▶ as $\lambda \rightarrow \infty$, slopes $\hat{\beta}_j^{(\lambda)} \rightarrow 0$ for $j = 1, \dots, p$.

▶ Use cross-validation to select λ .

Bayes use Prior to pick λ

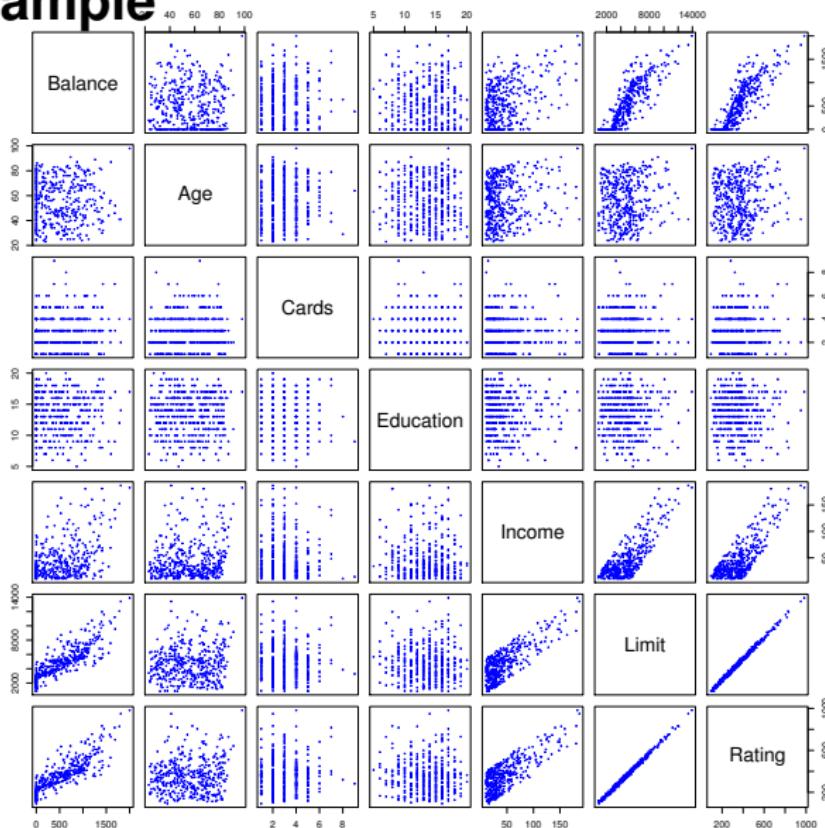
$\text{Var}(\hat{\beta}^{(\lambda)})$ vs $\text{Bias}(\hat{\beta}^{(\lambda)})$ Trade-off

Prior is similar to cross-validation

as $\lambda \rightarrow \infty$ we are Bias $\hat{\beta}^{(\lambda)}$ to be 0, so $\text{Var}(\hat{\beta}^{(\lambda)})$ is small

Credit example

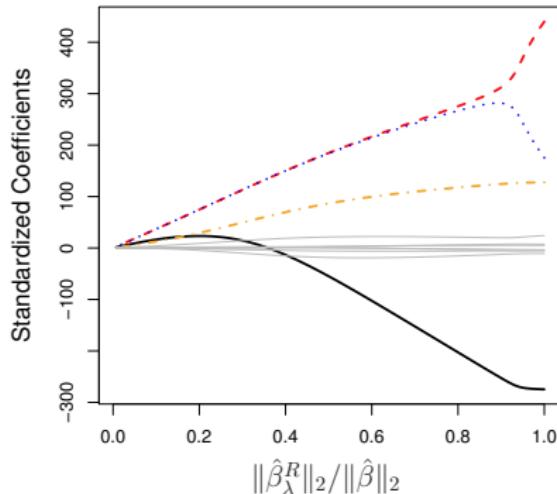
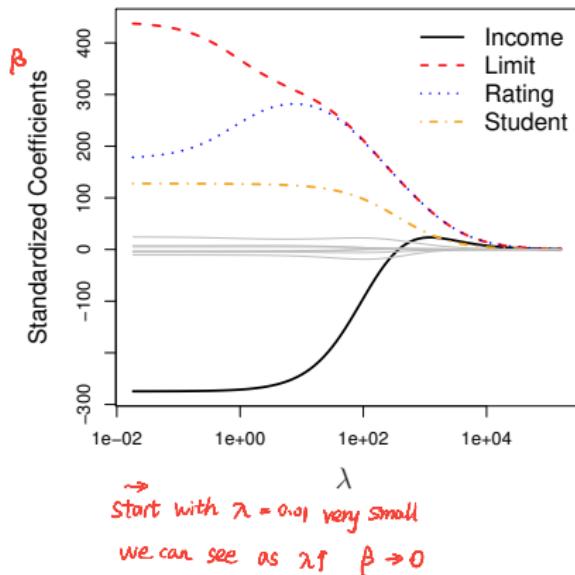
信用卡



Credit example

predict How much credit a client should get

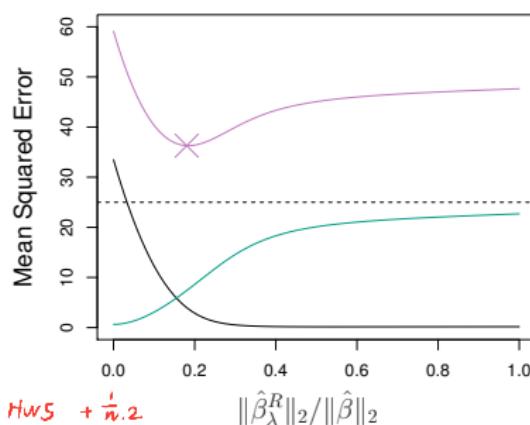
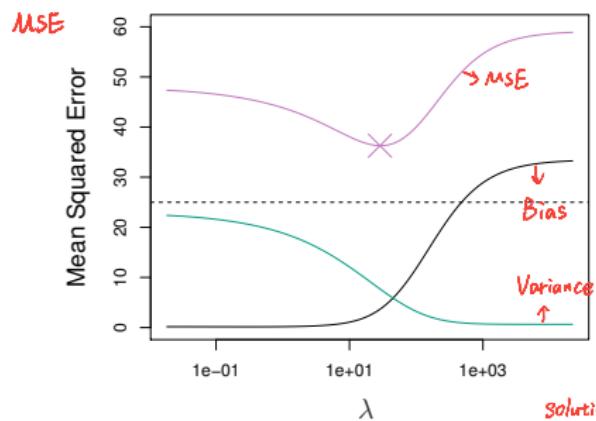
Ridge Regression solution



How does ridge regression work?

Consider model $y = X\beta + \epsilon$ where $\epsilon_{(i)}$ zero-mean/variance σ^2 .

MSE vs $\lambda \rightarrow$ Bias-Variance Tradeoff



- $E(\hat{\beta}^{(\lambda)}) = \beta - \lambda(X^T X + \lambda I)^{-1}\beta$. Not unbiased! $\text{if } \lambda=0 E=\beta \text{ unbiased}$
 $\lambda \uparrow \text{Bias} \uparrow$
- $\text{Var}(\hat{\beta}^{(\lambda)}) = \sigma^2(X^T X + \lambda I_p)^{-1}X^T X(X^T X + \lambda I_p)^{-1} \leq \text{Var}(\hat{\beta})$.
 $\lambda \uparrow \text{Var} \downarrow$

Scale-equivariant

For Linear Regression, Standardization doesn't matter but help numerically for large x calculation

Least squares coefficients are *scale-equivariant*.

Multiplying values of a predictor by a constant c simply results in multiplying its associated $\hat{\beta}_j$ by $1/c$.

What about for ridge regression?

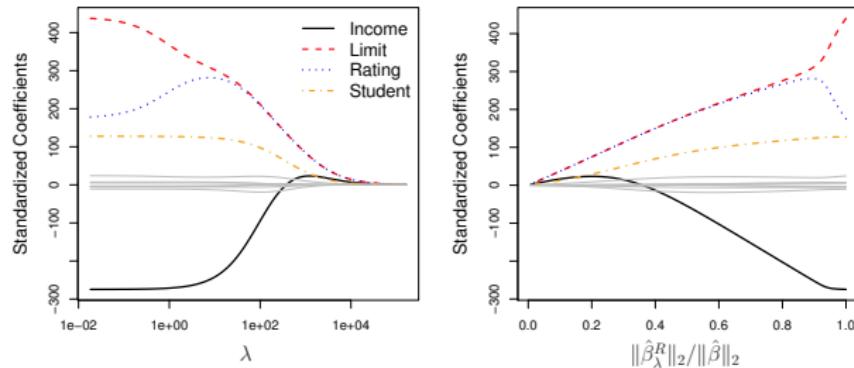
Not the case. Standardization is important
e.g. for KNN

$$(y - X\beta)^T(y - X\beta) + \lambda \beta^T \beta.$$

Because the shrinkage penalty depends on the magnitude of the β_j , scales now matter! So it's important we standardize the predictors.

continuous: to have mean=0 . std=1
categorical $x = -1, +1$

Credit example again



Note that for any choice of λ , $|\beta_j^{(\lambda)}| > 0$ for $j = 1, \dots, p$. This is fine for prediction, but not so good for interpretation.

↓
causality : A set of features can have accurate prediction of y
don't mean they cause y

Example

$$Y = \underbrace{3X_1 + \cdots + 3X_5}_{\text{relevant}} + 0X_6 + \cdots + 0X_{1000} + \epsilon$$

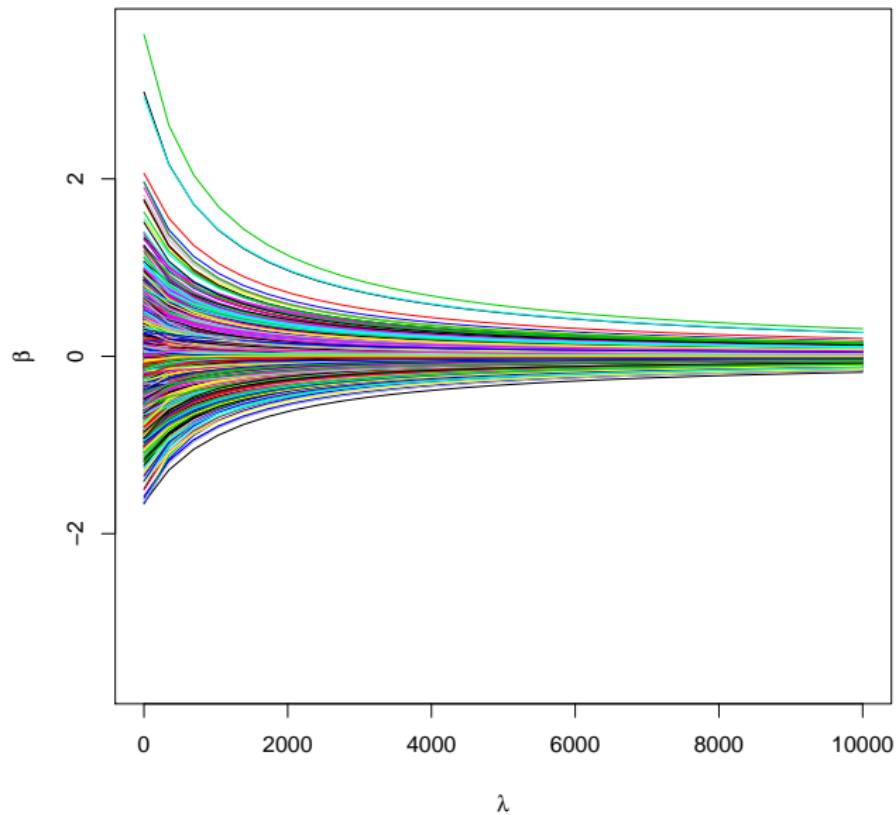
$$n = 100, p = 1,000.$$

So there are 1000 covariates but only 5 are relevant.

What does ridge regression do in this case?

A bad Predictor
instead we use Lasso

Ridge Regularization Paths



Lasso

Lasso (Least absolute shrinkage and selection operator) uses ℓ_1 penalty:

$$RSS(\beta) + \lambda \sum_{j=1}^p |\beta_j| ,$$

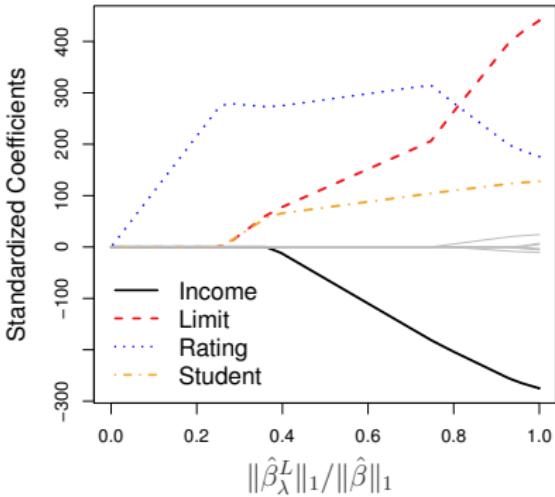
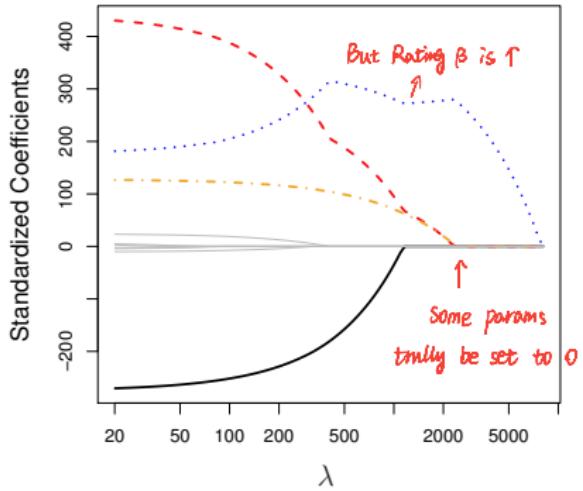
shrinkage penalty

No closed form solution for $\hat{\beta}^L$, but numerical optimization is “easy”

This modified penalized RSS has the effect of forcing some $\hat{\beta}_j^L$ to be exactly 0 for large λ .

① Gradient Descent
② threshold Gradient Descent

Credit example (with lasso)



An alternate view of lasso and ridge

Help understand geometric difference between ridge and Lasso why Lasso encourage solutions to be sparse
分枝的

An equivalent formulation of the ridge regression:

Constrained optimization

Residual Sum of Squares

$$RSS(\beta) = \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2$$

Constrained version ↓

subject to the constraint $\sum_{j=1}^p \beta_j^2 \leq s$ for some $s > 0$.

penalized version $\lambda \sum_{j=1}^p |\beta_j|$

Likewise, an equivalent formulation of lasso is to minimize $RSS(\beta)$

s.t. $\sum_{j=1}^p |\beta_j| \leq s$ for some $s > 0$.

such that

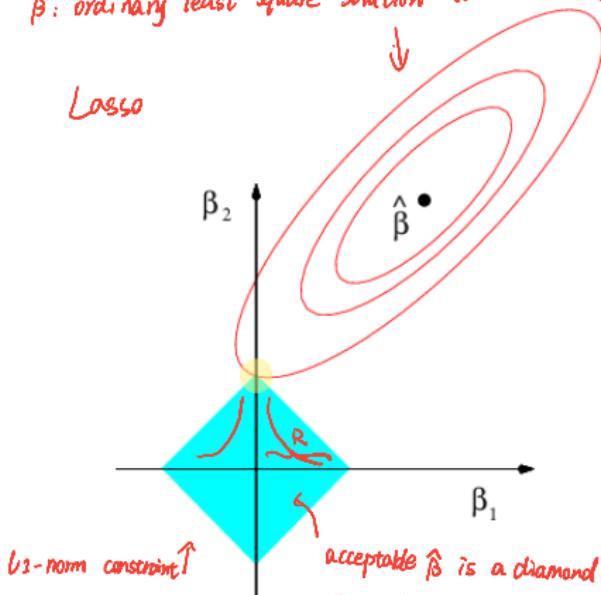
In each case, we can show that there's a 1-1 correspondence between s and λ .

Geometry of two methods

assume $\beta \in \mathbb{R}^2$

\diamond ellipse: least square objects

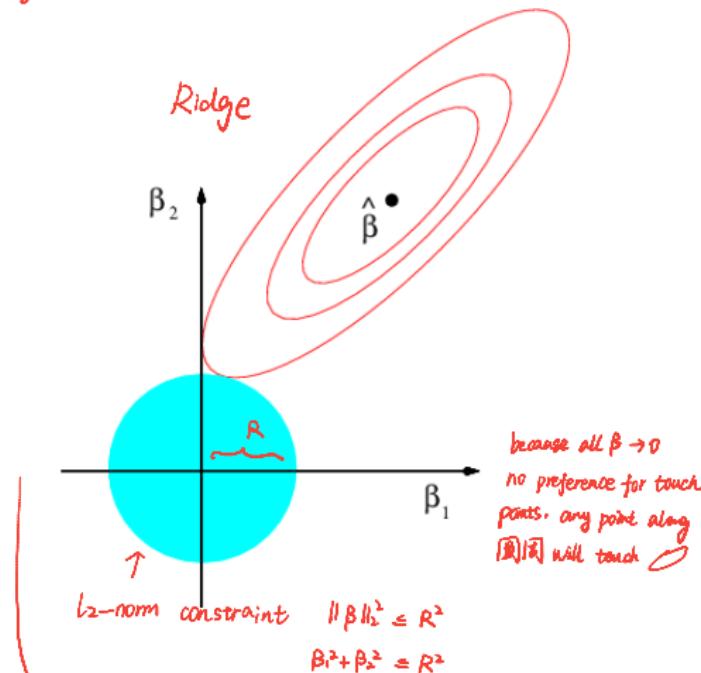
$\hat{\beta}$: ordinary least square solution with no constraint



Find β in the blue set which has smallest ellipse
some β_1 large some $\beta_2 \rightarrow 0$

Because \diamond is pointy, the first place \diamond and \square will touch may be one of points not the flat region. thus solution of Lasso is sparse

③ l_p norm $p \rightarrow \infty$
encourage $|\beta_1| = |\beta_2|$



Transforming predictors

Be careful when

How does transformation affect Lasso or Ridge?

- Important to scale for both *orthogonal transformation*
- But lasso is not orthogonally invariant $((X, \beta) \equiv (X\Gamma, \Gamma^\top \beta))$
 $\Gamma^\top \beta$ for a general Γ (gamma) *正交不變* *Γ*
• Destroys the sparsity pattern *orthogonal transform row of matrix X*
But for Ridge, no effect because l_2 -norm of a vector is invariant for orthogonal transformation

Summary of shrinkage methods

- General ℓ_q regularization criterion:

$$\hat{\beta} = \arg \min_{\beta} \left\{ \sum_{i=1}^n \underbrace{\left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2}_{\text{Loss function}} + \lambda \sum_{j=1}^p |\beta_j|^q \right\}$$

*↑
L_q norm regularizer*

What would ℓ_0 be doing? *non-convex* = it counts # non-zero coefficients

- Idea can apply to many methods, including classification. For example, logistic regression, modify negative log-likelihood problem to:

By change Loss function

$$\frac{1}{n} \sum_{i=1}^n \left[\log(1 + \exp(-x_i^T \beta)) - y_i x_i^T \beta \right] + \frac{\lambda}{2} \sum_{j=1}^p \beta_j^2$$

- Compared to ℓ_1 (lasso), ℓ_2 (ridge) produces models that are less interpretable because there will be a large number of coefficients.

Comparing lasso and ridge

An equivalent formulation of the ridge regression optimization problem is that we want to minimize:

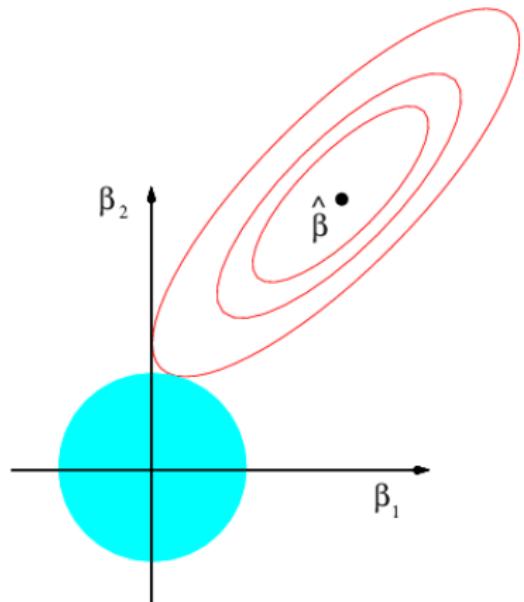
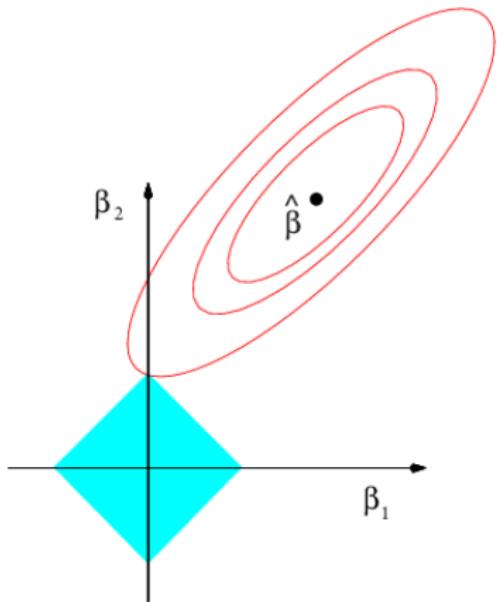
$$RSS(\beta) = \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2$$

subject to the constraint $\sum_{j=1}^p \beta_j^2 \leq s$ for some $s > 0$.

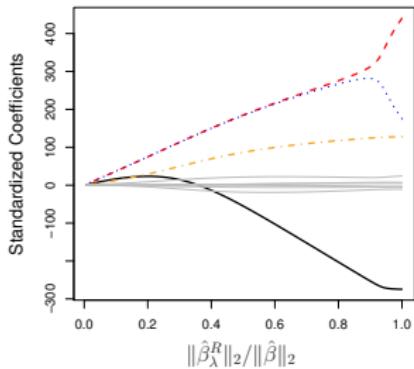
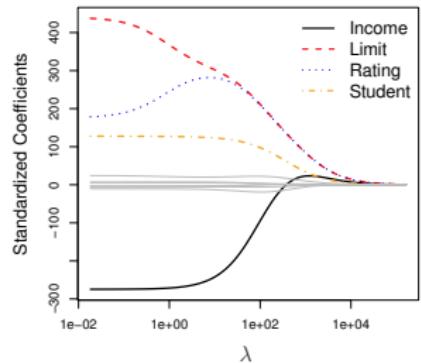
Likewise, an equivalent formulation of lasso is to minimize $RSS(\beta)$
s.t. $\sum_{j=1}^p |\beta_j| \leq s$ for some $s > 0$.

In each case, we can show that there's a 1-1 correspondence between s and λ .

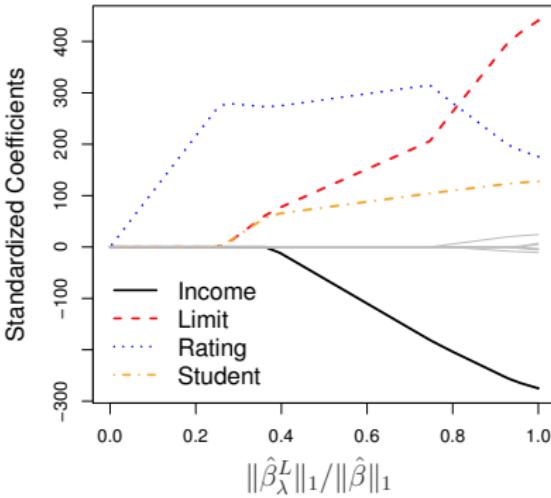
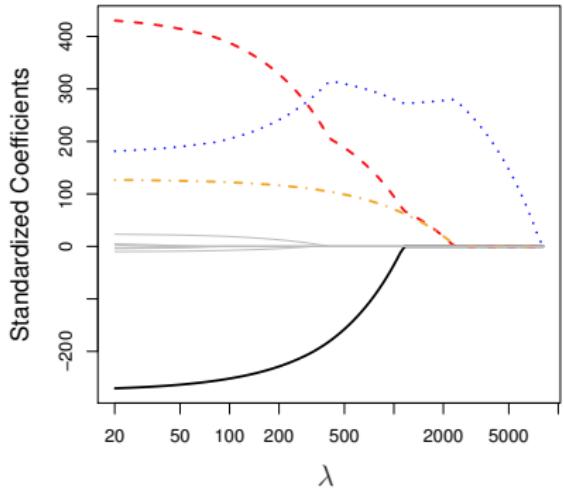
Comparing lasso and ridge



Credit example again



Credit example (with lasso)



How to choose λ

- Cross-validation is straight-forward:
 - ▶ Set up a grid of λ values.
 - ▶ Estimate test error using 10-fold cross-validation for each λ .
 - ▶ Choose λ with lowest cross-validation error.

Hw5 Boosting test error 5%.

Comparing lasso and ridge

Solution to HW5

Consider the special case where we have orthonormal predictors (and no intercept). That is to say, that $X^T X = I$. Then, the lasso coefficients are

where $\hat{\beta}_j$ is least square solution

$$\hat{\beta}_j^L = \begin{cases} \hat{\beta}_j - \lambda & \text{if } \hat{\beta}_j > \lambda \\ 0 & \text{if } |\hat{\beta}_j| \leq \lambda \\ \hat{\beta}_j + \lambda & \text{if } \hat{\beta}_j < -\lambda \end{cases}$$

(Additive Shrinkage) ↓ Minus ↑ Identity

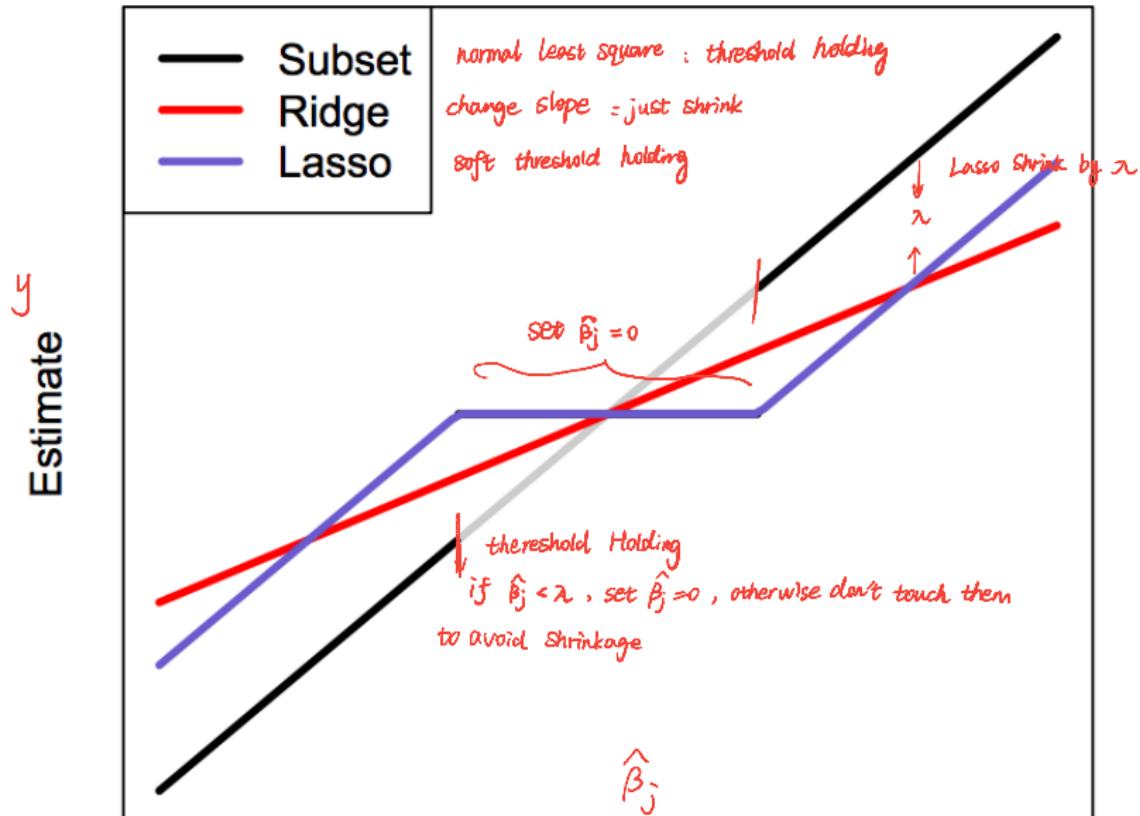
And the ridge coefficients are:

(Multiplicative Shrinkage)

$$\hat{\beta}_j^{(\lambda)} = \frac{\hat{\beta}_j}{1 + \lambda}.$$

Problem with Lasso : if true β_j is very large , $\hat{\beta}_j^L = \hat{\beta}_j - \lambda$ might penalty too much by λ . e.g. $\hat{\beta}_j = 10 \rightarrow \hat{\beta}_j^L = 1$ if you know $\beta_j = 10$, you should be very confident that $\hat{\beta}_j = 10$. So we also use Bayes method to do variable selection to take this phenomenon into account : If you are very confident with $\hat{\beta}_j$, you won't shrink it by λ , similar to Ridge

Comparing lasso and ridge

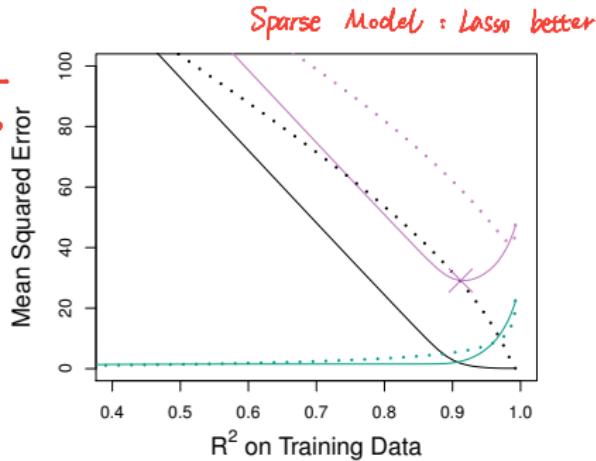
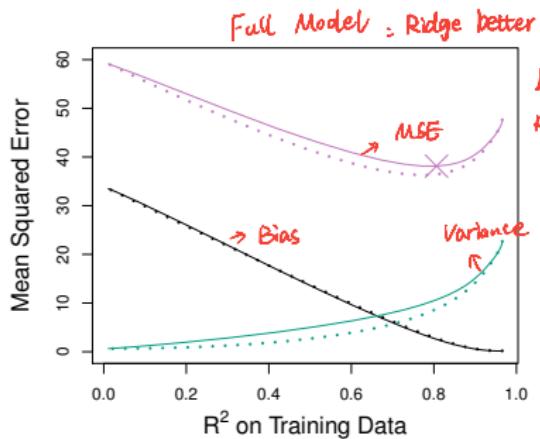


Comparing lasso and ridge

- Lasso yields simpler (sparser) models.
- Specifically, lasso can select at most n features.
- Lasso and ridge deal with multicollinearity differently.
 - just shrink all features down
↓ pick one of strongest features
- Actual performance (MSE) depends on underlying truth.
Use Cross validation or Test / Train Split to assess performance

Comparing lasso and ridge

Simulated data: 45 predictors.



- Lasso (solid) vs. ridge (dotted) lines
- Bias² (black), variance (green), test MSE (purple)
- Generating model:

full model ▶ left: $y = \beta_0 + \beta_1 x_1 + \dots + \beta_{45} x_{45} + \epsilon$

sparse model (reduced) ▶ right: $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$ only use 2 predictors

Often using an ℓ_2 penalty encourages better prediction performance.
For that reason ℓ_1 and ℓ_2 penalties are combined to form **elastic net**

Lasso Ridge 岭回归 弹性网络

$$2 \text{ hyper params } \lambda, \mu \quad \arg \min_{\beta} \mathcal{L}(\beta) + \lambda \|\beta\|_1 + \mu \|\beta\|_2^2$$

\uparrow \uparrow

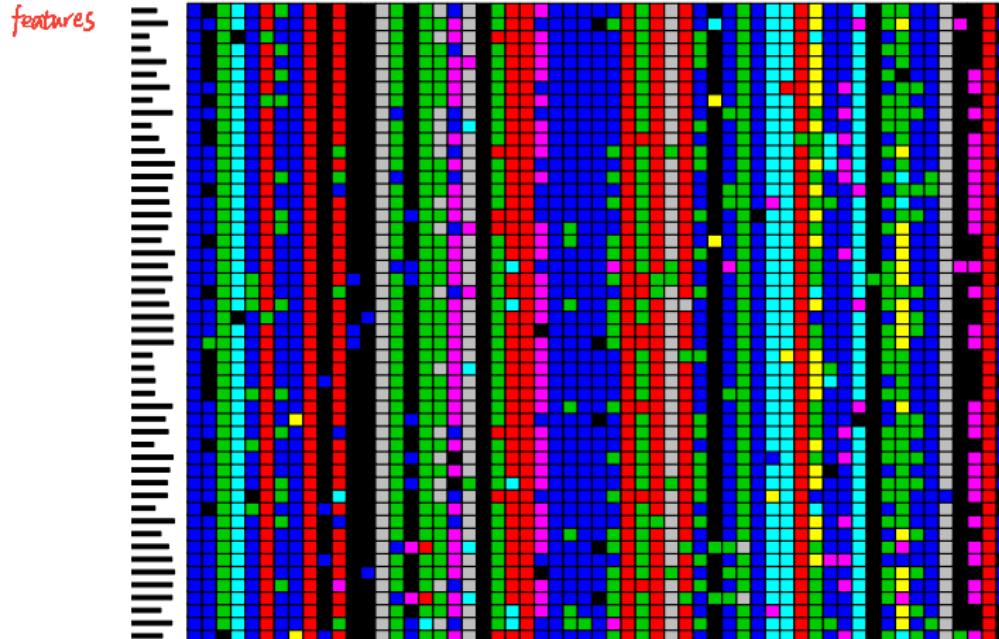
$\ell_1\text{-norm}$ $\ell_2\text{-norm}$

$\left\{ \begin{array}{l} \text{hyper params: params setted before training} \\ \text{params: } \beta \text{ params obtained after training} \end{array} \right.$

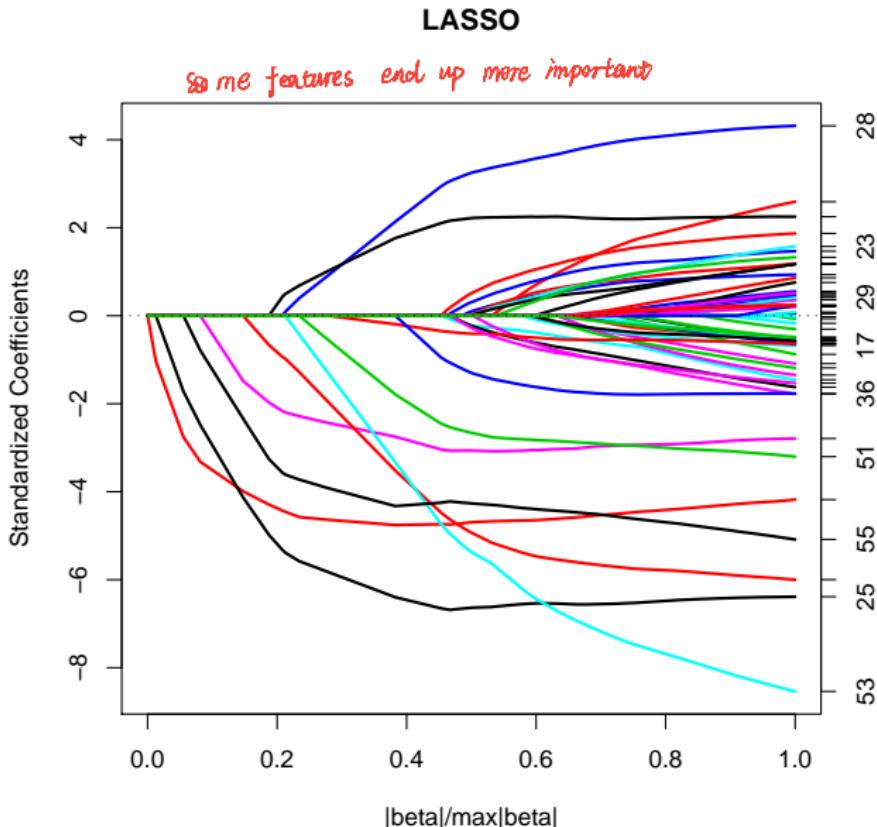
pick optimal λ and μ using cross-validation by testing combination of λ and μ

HIV example

- Y is resistance to HIV drug.
- X_j = amino acid in position j of the virus.
- $p = 99$, $n \approx 100$. *samples*

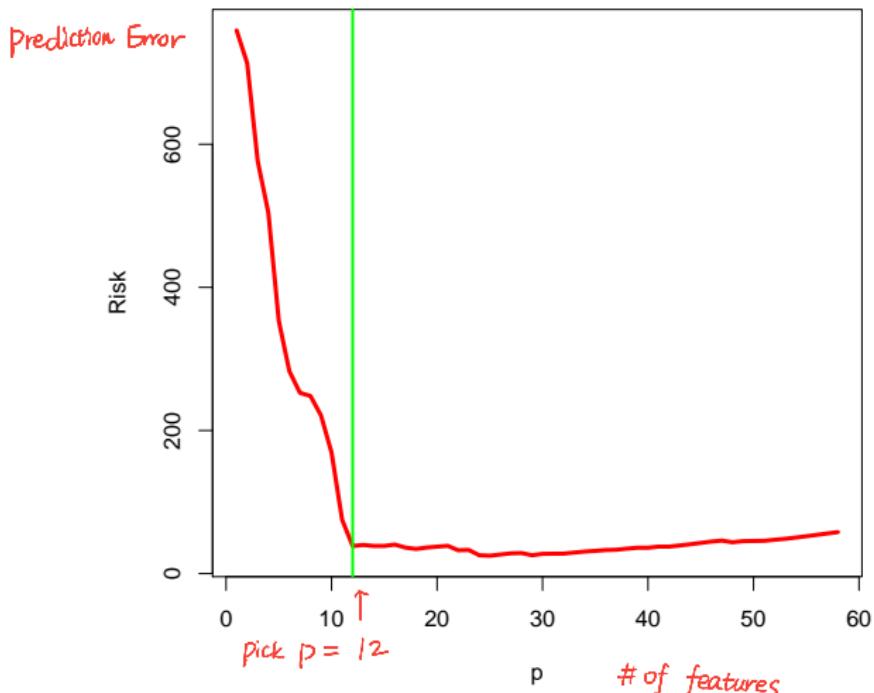


HIV example

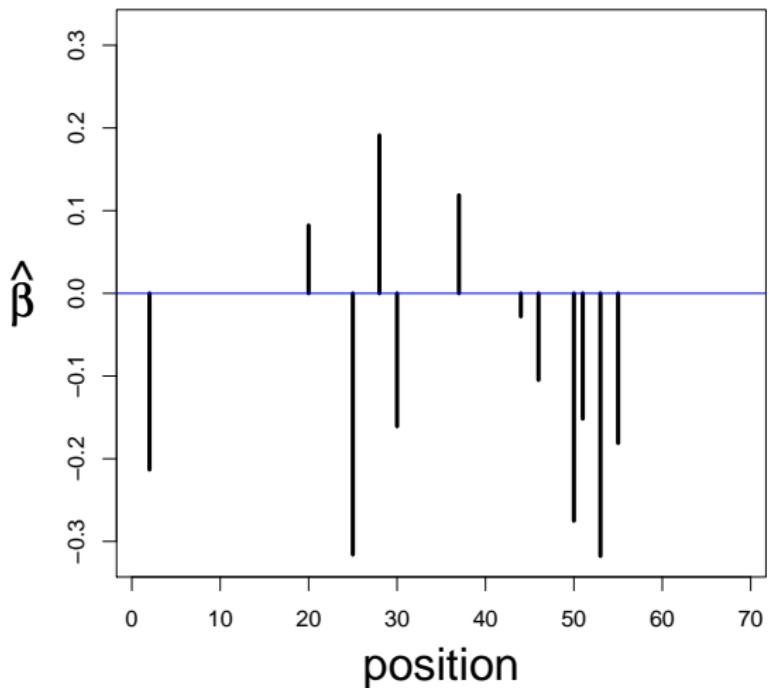


Selecting λ

We choose the sparsity level by estimating prediction error.
(# of features)



HIV example



Conclusion

- scale matters standardization is important
- Ridge regression has closed form
- Penalizing by ℓ_1 norm leads to sparsity (and shrinkage)
feature 变少 β 变小
- If the true model only involves a few predictors, lasso may be better than ridge

Q: When to use ridge is better than Lasso?

$$\begin{array}{ccc} \downarrow & & \downarrow \\ l_2 & & l_2 \end{array}$$

{
Lasso : help robustness
Ridge : help robustness, help protect from noise, better prediction accuracy
Elastic net : combine both strengths