

LECTURE 11: REGULARIZATION

STAT 598Z: INTRODUCTION TO COMPUTING FOR STATISTICS

Vinayak Rao

Department of Statistics, Purdue University

February 15, 2017

ORDINARY LEAST SQUARES

Consider linear regression:

$$y = \mathbf{x}^T \mathbf{w} + \epsilon$$

A diagram illustrating the linear regression equation $y = \mathbf{x}^T \mathbf{w} + \epsilon$. It shows three colored rectangles: a blue square representing y , a cyan rectangle representing x^T , and an orange rectangle representing w . The equation is written above the rectangles: $y = x^T w$. The blue square is on the left, followed by an equals sign, then the cyan rectangle, then the orange rectangle. The labels y , x^T , and w are placed above their respective rectangles.

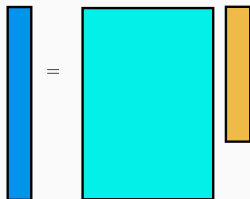
ORDINARY LEAST SQUARES

Consider linear regression:

$$y = \mathbf{x}^\top \mathbf{w} + \epsilon$$

In vector notation:

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \epsilon, \quad \mathbf{y} \in \mathbb{R}^n, \mathbf{X} \in \mathbb{R}^{n \times p}$$



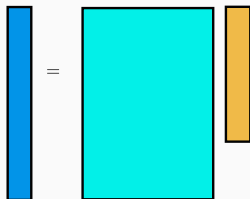
ORDINARY LEAST SQUARES

Consider linear regression:

$$y = \mathbf{x}^\top \mathbf{w} + \epsilon$$

In vector notation:

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \epsilon, \quad \mathbf{y} \in \mathbb{R}^n, \mathbf{X} \in \mathbb{R}^{n \times p}$$



$$\hat{\mathbf{w}} = \arg \min \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|^2 = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{w})^2$$

ORDINARY LEAST SQUARES

Problem:

$$\hat{\mathbf{w}} = \arg \min \|\mathbf{y} - \mathbf{X}^T \mathbf{w}\|^2 = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \mathbf{w})^2$$

Solution:

$$\hat{\mathbf{w}} = (\mathbf{X}\mathbf{X}^T)^{-1}\mathbf{X}\mathbf{y} \quad (\text{correlation in 1-d})$$

ORDINARY LEAST SQUARES

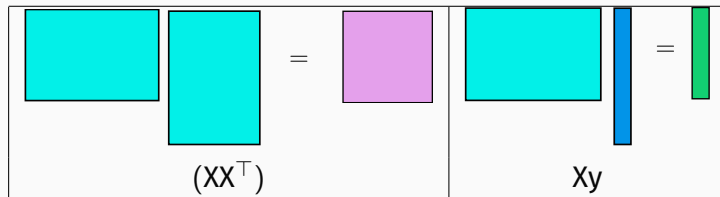
Problem:

$$\hat{\mathbf{w}} = \arg \min \|\mathbf{y} - \mathbf{X}^T \mathbf{w}\|^2 = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \mathbf{w})^2$$

Solution:

$$\hat{\mathbf{w}} = (\mathbf{X}\mathbf{X}^T)^{-1}\mathbf{X}\mathbf{y}$$

(correlation in 1-d)



How to do this in R (without using `lm`)?

- Do not invert with `solve` and multiply!

ORDINARY LEAST SQUARES

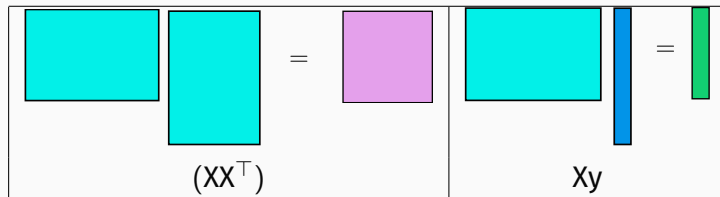
Problem:

$$\hat{\mathbf{w}} = \arg \min \|\mathbf{y} - \mathbf{X}^T \mathbf{w}\|^2 = \arg \min \sum_{i=1}^n (y_i - \mathbf{x}_i^T \mathbf{w})^2$$

Solution:

$$\hat{\mathbf{w}} = (\mathbf{X}\mathbf{X}^T)^{-1}\mathbf{X}\mathbf{y}$$

(correlation in 1-d)



How to do this in R (without using `lm`)?

- Do not invert with `solve` and multiply!
- Directly solve $(\mathbf{X}\mathbf{X}^T)\hat{\mathbf{w}} = \mathbf{X}\mathbf{y}$

PREDICTION ERROR

$\hat{\mathbf{w}}$ is an unbiased estimate of the true \mathbf{w}

For a test vector \mathbf{x}^{test} we predict $\mathbf{w}^\top \mathbf{x}^{test}$.

(Squared) prediction error: $PE = \frac{1}{k} \sum_{i=1}^k (y_i^{test} - \mathbf{w}^\top \mathbf{x}_i^{test})^2$

$\hat{\mathbf{w}}$ is an unbiased estimate of the true \mathbf{w}

For a test vector \mathbf{x}^{test} we predict $\mathbf{w}^\top \mathbf{x}^{test}$.

(Squared) prediction error: $PE = \frac{1}{k} \sum_{i=1}^k (y_i^{test} - \mathbf{w}^\top \mathbf{x}_i^{test})^2$

Can show:

- PE is has mean 0
- variance grows with number of features (p)

$\hat{\mathbf{w}}$ is an unbiased estimate of the true \mathbf{w}

For a test vector \mathbf{x}^{test} we predict $\mathbf{w}^\top \mathbf{x}^{test}$.

(Squared) prediction error: $PE = \frac{1}{k} \sum_{i=1}^k (y_i^{test} - \mathbf{w}^\top \mathbf{x}_i^{test})^2$

Can show:

- PE is has mean 0
- variance grows with number of features (p)

What if $p > n$?

- $\mathbf{X}\mathbf{X}^\top$ is singular

$p > n$:

- Cannot invert $\mathbf{X}\mathbf{X}^\top$

$p > n$:

- Cannot invert $\mathbf{X}\mathbf{X}^\top$
- We *can* invert if we add a small λ to the diagonal

$$\hat{\mathbf{w}}_\lambda = (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I})^{-1} \mathbf{X}\mathbf{y} \quad (\mathbf{I} \text{ is the identity matrix})$$

$p > n$:

- Cannot invert $\mathbf{X}\mathbf{X}^\top$
- We *can* invert if we add a small λ to the diagonal

$$\hat{\mathbf{w}}_\lambda = (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I})^{-1} \mathbf{X}\mathbf{y} \quad (\mathbf{I} \text{ is the identity matrix})$$

Introducing λ makes problem well-posed, but introduces bias

$p > n$:

- Cannot invert $\mathbf{X}\mathbf{X}^\top$
- We *can* invert if we add a small λ to the diagonal

$$\hat{\mathbf{w}}_\lambda = (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I})^{-1} \mathbf{X}\mathbf{y} \quad (\mathbf{I} \text{ is the identity matrix})$$

Introducing λ makes problem well-posed, but introduces bias

- $\lambda = 0$ recovers OLS
- Larger λ causes larger bias

$p > n$:

- Cannot invert $\mathbf{X}\mathbf{X}^\top$
- We *can* invert if we add a small λ to the diagonal

$$\hat{\mathbf{w}}_\lambda = (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I})^{-1} \mathbf{X}\mathbf{y} \quad (\mathbf{I} \text{ is the identity matrix})$$

Introducing λ makes problem well-posed, but introduces bias

- $\lambda = 0$ recovers OLS
- Larger λ causes larger bias
- $\lambda = \infty$?

$p > n$:

- Cannot invert $\mathbf{X}\mathbf{X}^\top$
- We *can* invert if we add a small λ to the diagonal

$$\hat{\mathbf{w}}_\lambda = (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I})^{-1} \mathbf{X}\mathbf{y} \quad (\mathbf{I} \text{ is the identity matrix})$$

Introducing λ makes problem well-posed, but introduces bias

- $\lambda = 0$ recovers OLS
- Larger λ causes larger bias
- $\lambda = \infty$? No variance!

$p > n$:

- Cannot invert $\mathbf{X}\mathbf{X}^\top$
- We *can* invert if we add a small λ to the diagonal

$$\hat{\mathbf{w}}_\lambda = (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I})^{-1} \mathbf{X}\mathbf{y} \quad (\mathbf{I} \text{ is the identity matrix})$$

Introducing λ makes problem well-posed, but introduces bias

- $\lambda = 0$ recovers OLS
- Larger λ causes larger bias
- $\lambda = \infty$? No variance!

λ trades-off bias and variance

Maybe a nonzero λ is actually good?

RIDGE REGRESSION (A.K.A. TIKHONOV REGULARIZATION)

Recall $\hat{\mathbf{w}} = (\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{X}\mathbf{y}$ solves $\hat{\mathbf{w}} = \arg \min \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|^2$

RIDGE REGRESSION (A.K.A. TIKHONOV REGULARIZATION)

Recall $\hat{\mathbf{w}} = (\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{X}\mathbf{y}$ solves $\hat{\mathbf{w}} = \arg \min \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|^2$

$\hat{\mathbf{w}}_\lambda = (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I})^{-1}\mathbf{X}\mathbf{y}$ solves

$$\hat{\mathbf{w}}_\lambda = \operatorname{argmin} \mathcal{L}_\lambda(\mathbf{w}) := \operatorname{argmin} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2$$

RIDGE REGRESSION (A.K.A. TIKHONOV REGULARIZATION)

Recall $\hat{\mathbf{w}} = (\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{X}\mathbf{y}$ solves $\hat{\mathbf{w}} = \arg \min \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|^2$

$\hat{\mathbf{w}}_\lambda = (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I})^{-1}\mathbf{X}\mathbf{y}$ solves

$$\hat{\mathbf{w}}_\lambda = \operatorname{argmin} \mathcal{L}_\lambda(\mathbf{w}) := \operatorname{argmin} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2$$

$\|\mathbf{w}\|_2^2 = \sum_{i=1}^p w_i^2$ is the squared ℓ_2 -norm

$\lambda \|\mathbf{w}\|_2$ is the *shrinkage penalty*.

RIDGE REGRESSION (A.K.A. TIKHONOV REGULARIZATION)

Recall $\hat{\mathbf{w}} = (\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{X}\mathbf{y}$ solves $\hat{\mathbf{w}} = \arg \min \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|^2$

$\hat{\mathbf{w}}_\lambda = (\mathbf{X}\mathbf{X}^\top + \lambda I)^{-1}\mathbf{X}\mathbf{y}$ solves

$$\hat{\mathbf{w}}_\lambda = \operatorname{argmin} \mathcal{L}_\lambda(\mathbf{w}) := \operatorname{argmin} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2$$

$\|\mathbf{w}\|_2^2 = \sum_{i=1}^p w_i^2$ is the squared ℓ_2 -norm

$\lambda \|\mathbf{w}\|_2$ is the *shrinkage penalty*.

Favours \mathbf{w} 's with smaller components

RIDGE REGRESSION (A.K.A. TIKHONOV REGULARIZATION)

Recall $\hat{\mathbf{w}} = (\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{X}\mathbf{y}$ solves $\hat{\mathbf{w}} = \arg \min \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|^2$

$\hat{\mathbf{w}}_\lambda = (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I})^{-1}\mathbf{X}\mathbf{y}$ solves

$$\hat{\mathbf{w}}_\lambda = \operatorname{argmin} \mathcal{L}_\lambda(\mathbf{w}) := \operatorname{argmin} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2$$

$\|\mathbf{w}\|_2^2 = \sum_{i=1}^p w_i^2$ is the squared ℓ_2 -norm

$\lambda \|\mathbf{w}\|_2$ is the *shrinkage penalty*.

Favours \mathbf{w} 's with smaller components

λ trades off small training error with 'simple' solutions

RIDGE REGRESSION (A.K.A. TIKHONOV REGULARIZATION)

Recall $\hat{\mathbf{w}} = (\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{X}\mathbf{y}$ solves $\hat{\mathbf{w}} = \arg \min \|\mathbf{y} - \mathbf{X}^\top \mathbf{w}\|^2$

$\hat{\mathbf{w}}_\lambda = (\mathbf{X}\mathbf{X}^\top + \lambda I)^{-1}\mathbf{X}\mathbf{y}$ solves

$$\hat{\mathbf{w}}_\lambda = \operatorname{argmin} \mathcal{L}_\lambda(\mathbf{w}) := \operatorname{argmin} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2$$

$\|\mathbf{w}\|_2^2 = \sum_{i=1}^p w_i^2$ is the squared ℓ_2 -norm

$\lambda \|\mathbf{w}\|_2$ is the *shrinkage penalty*.

Favours \mathbf{w} 's with smaller components

λ trades off small training error with 'simple' solutions

ℓ_2 /ridge/Tikhonov regularization

Simple modification of the least-squares solution:

$$\hat{\mathbf{w}}_{\lambda} = (\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^{\top} \mathbf{y}$$

RIDGE REGRESSION (SOLUTION)

Simple modification of the least-squares solution:

$$\hat{\mathbf{w}}_{\lambda} = (\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^{\top} \mathbf{y}$$

In the 1-dimensional case,

$$\hat{w}_{\lambda} = (\mathbf{x}^{\top} \mathbf{x} + \lambda \mathbf{I}_p)^{-1} \mathbf{x}^{\top} \mathbf{y}$$

RIDGE REGRESSION (SOLUTION)

Simple modification of the least-squares solution:

$$\hat{\mathbf{w}}_{\lambda} = (\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^{\top} \mathbf{y}$$

In the 1-dimensional case,

$$\begin{aligned} \hat{w}_{\lambda} &= (\mathbf{x}^{\top} \mathbf{x} + \lambda \mathbf{I}_p)^{-1} \mathbf{x}^{\top} \mathbf{y} \\ &= \frac{\mathbf{x}^{\top} \mathbf{x}}{(\mathbf{x}^{\top} \mathbf{x} + \lambda \mathbf{I}_p)} \frac{\mathbf{x}^{\top} \mathbf{y}}{\mathbf{x}^{\top} \mathbf{x}} \end{aligned}$$

RIDGE REGRESSION (SOLUTION)

Simple modification of the least-squares solution:

$$\hat{\mathbf{w}}_\lambda = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{y}$$

In the 1-dimensional case,

$$\begin{aligned}\hat{w}_\lambda &= (\mathbf{x}^\top \mathbf{x} + \lambda \mathbf{I}_p)^{-1} \mathbf{x}^\top \mathbf{y} \\ &= \frac{\mathbf{x}^\top \mathbf{x}}{(\mathbf{x}^\top \mathbf{x} + \lambda \mathbf{I}_p)} \frac{\mathbf{x}^\top \mathbf{y}}{\mathbf{x}^\top \mathbf{x}} \\ &= c \hat{w} \quad (c < 1)\end{aligned}$$

RIDGE REGRESSION (SOLUTION)

Simple modification of the least-squares solution:

$$\hat{\mathbf{w}}_\lambda = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{y}$$

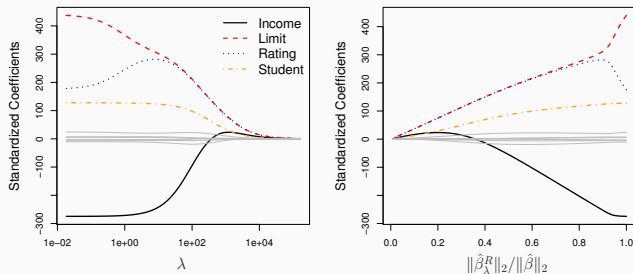
In the 1-dimensional case,

$$\begin{aligned}\hat{w}_\lambda &= (\mathbf{x}^\top \mathbf{x} + \lambda \mathbf{I}_p)^{-1} \mathbf{x}^\top \mathbf{y} \\ &= \frac{\mathbf{x}^\top \mathbf{x}}{(\mathbf{x}^\top \mathbf{x} + \lambda \mathbf{I}_p)} \frac{\mathbf{x}^\top \mathbf{y}}{\mathbf{x}^\top \mathbf{x}} \\ &= c \hat{w} \quad (c < 1)\end{aligned}$$

Shrinks least-squares solution.

RIDGE REGRESSION

Credit data set (average credit card debt)



James, Witten, Hastie and Tibshirani

HOW DO WE CHOOSE λ ?

Cross-validation:

HOW DO WE CHOOSE λ ?

Cross-validation:

- Pick a set of λ 's
- For k th fold of cross-validation:

HOW DO WE CHOOSE λ ?

Cross-validation:

- Pick a set of λ 's
- For k th fold of cross-validation:
 - For each λ :
 - Solve the regularized least squares problem on training data.
 - Evaluate estimated \mathbf{w} on held-out data (call this $PE_{\lambda,k}$).

HOW DO WE CHOOSE λ ?

Cross-validation:

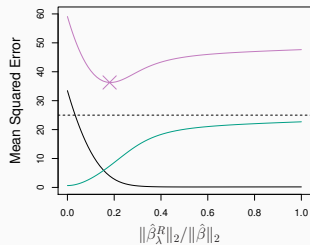
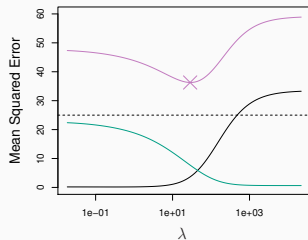
- Pick a set of λ 's
- For k th fold of cross-validation:
 - For each λ :
 - Solve the regularized least squares problem on training data.
 - Evaluate estimated \mathbf{w} on held-out data (call this $PE_{\lambda,k}$).
- Pick $\hat{\lambda} = \operatorname{argmin} \operatorname{mean}(PE_{\lambda})$
or $(\operatorname{argmin} (\operatorname{mean}(PE_{\lambda}) + \operatorname{stderr}(PE_{\lambda})))$

HOW DO WE CHOOSE λ ?

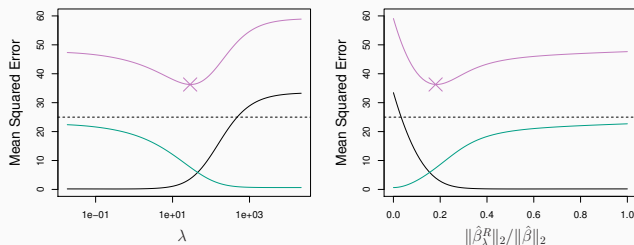
Cross-validation:

- Pick a set of λ 's
- For k th fold of cross-validation:
 - For each λ :
 - Solve the regularized least squares problem on training data.
 - Evaluate estimated \mathbf{w} on held-out data (call this $PE_{\lambda,k}$).
- Pick $\hat{\lambda} = \operatorname{argmin} \operatorname{mean}(PE_{\lambda})$
or $(\operatorname{argmin} (\operatorname{mean}(PE_{\lambda}) + \operatorname{stderr}(PE_{\lambda})))$
- Having chosen $\hat{\lambda}$ solve regularized least square on all data

DOES THIS WORK?

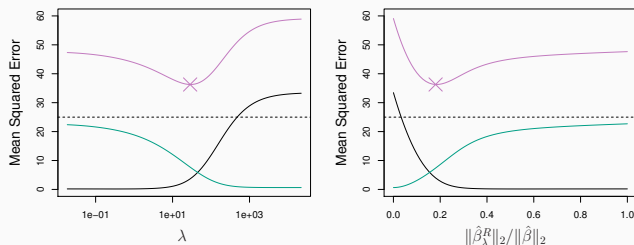


DOES THIS WORK?



Ridge regression improves performance by reducing variance

DOES THIS WORK?



Ridge regression improves performance by reducing variance

- does not perform feature selection
- just shrinks components of \mathbf{w} towards 0

For the former: Lasso

$\operatorname{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2$ is equivalent to

$\operatorname{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 \quad \text{s.t.} \quad \|\mathbf{w}\|_2^2 \leq \gamma$

(Note: γ will depend on data)

REGULARIZATION AS CONSTRAINED OPTIMIZATION

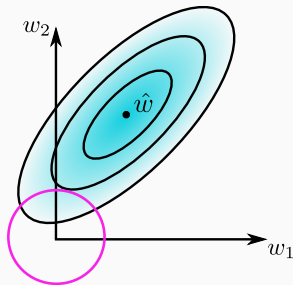
$\operatorname{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2$ is equivalent to

$\operatorname{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 \quad \text{s.t.} \quad \|\mathbf{w}\|_2^2 \leq \gamma$

(Note: γ will depend on data)

First problem: regularized optimization

Second problem: constrained optimization



REGULARIZATION AS CONSTRAINED OPTIMIZATION

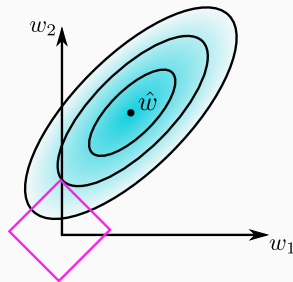
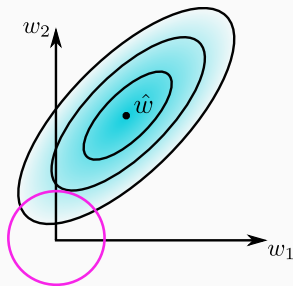
$\operatorname{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2$ is equivalent to

$\operatorname{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 \quad \text{s.t.} \quad \|\mathbf{w}\|_2^2 \leq \gamma$

(Note: γ will depend on data)

First problem: regularized optimization

Second problem: constrained optimization



$$\operatorname{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 \quad \text{s.t. } \|\mathbf{w}\|_1 \leq \gamma$$

$\|\mathbf{w}\|_1 = \sum_{i=1}^p |w_i|$ is the ℓ_1 -norm.

$$\operatorname{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 \quad \text{s.t. } \|\mathbf{w}\|_1 \leq \gamma$$

$\|\mathbf{w}\|_1 = \sum_{i=1}^p |w_i|$ is the ℓ_1 -norm.

Lasso: least absolute shrinkage and selection operator.

$$\hat{\mathbf{w}} = \operatorname{argmin} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_1$$

$$\operatorname{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 \quad \text{s.t. } \|\mathbf{w}\|_1 \leq \gamma$$

$\|\mathbf{w}\|_1 = \sum_{i=1}^p |w_i|$ is the ℓ_1 -norm.

Lasso: least absolute shrinkage and selection operator.

$$\hat{\mathbf{w}} = \operatorname{argmin} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_1$$

- Penalizes small w_i 's more than ridge regression.
- Tolerates larger w_i 's more than ridge regression.

$$\operatorname{argmin}(\mathbf{y} - \mathbf{X}^\top \mathbf{w})^2 \quad \text{s.t. } \|\mathbf{w}\|_1 \leq \gamma$$

$\|\mathbf{w}\|_1 = \sum_{i=1}^p |w_i|$ is the ℓ_1 -norm.

Lasso: least absolute shrinkage and selection operator.

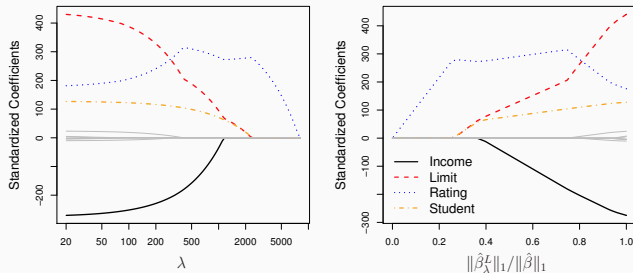
$$\hat{\mathbf{w}} = \operatorname{argmin} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_1$$

- Penalizes small w_i 's more than ridge regression.
- Tolerates larger w_i 's more than ridge regression.

Result:

- $\hat{\mathbf{w}}_{\text{LASSO}}$ has some components *exactly* equal to zero.
- Performs feature selection.

Credit data set (average credit card debt)



James, Witten, Hastic and Tibshirani