

Ridge Regression Recap

CS 229: Machine Learning

Emily Fox

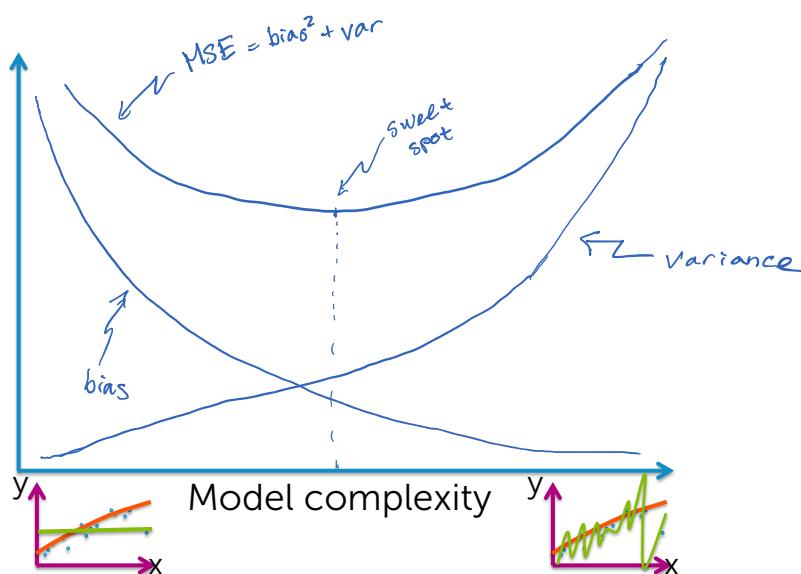
Stanford University

January 22, 2024

©2024 Emily Fox

1

Bias-variance tradeoff



©2024 Emily Fox

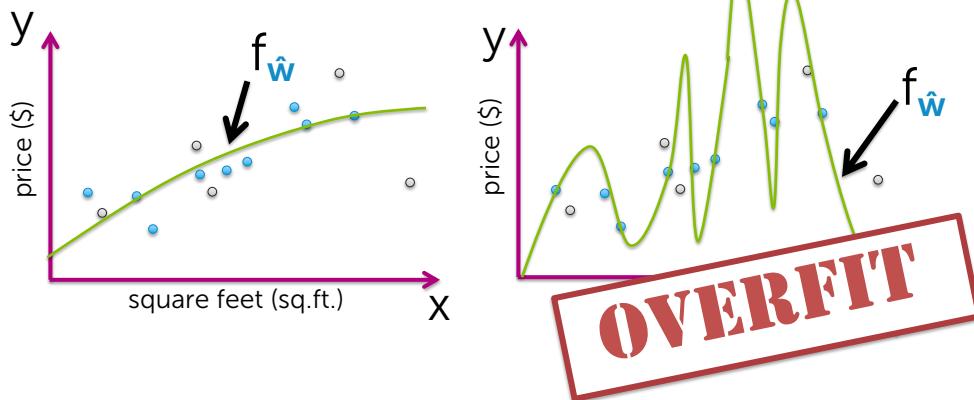
CS 229: Machine Learning

2

1

Flexibility of complex models (large D for linear regression)

$$y_i = \sum_{j=0}^D w_j h_j(x_i) + \epsilon_i$$



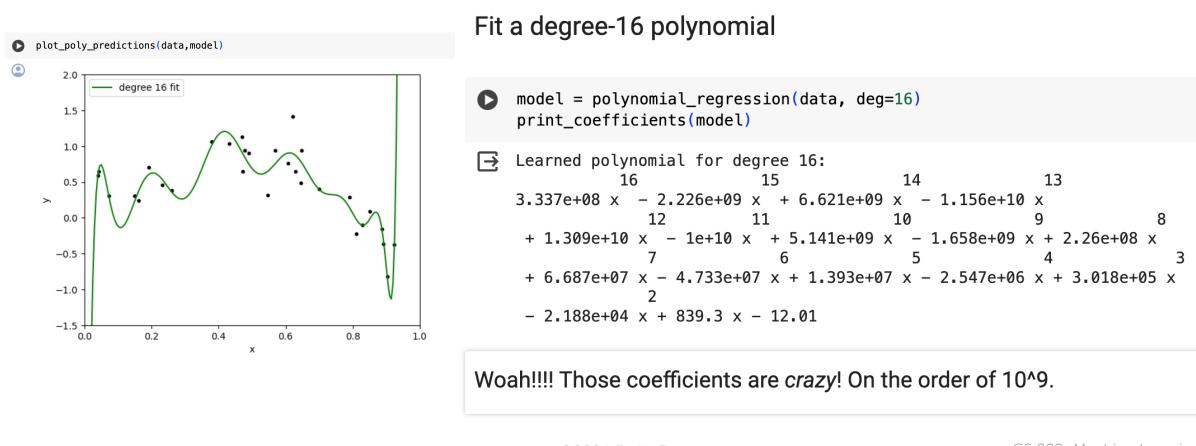
©2024 Emily Fox

CS 229: Machine Learning

3

Symptom of overfitting

Often, overfitting associated with very large estimated \hat{w}



©2024 Emily Fox

CS 229: Machine Learning

4

Consider specific total cost

Total cost =

$$\underbrace{\text{measure of fit}}_{\text{RSS}(\mathbf{w})} + \underbrace{\text{measure of magnitude of coefficients}}_{\|\mathbf{w}\|_2^2}$$

©2024 Emily Fox

CS 229: Machine Learning

5

Consider resulting objective

What if $\hat{\mathbf{w}}$ selected to minimize

$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$$

↑ tuning parameter = balance of fit and magnitude

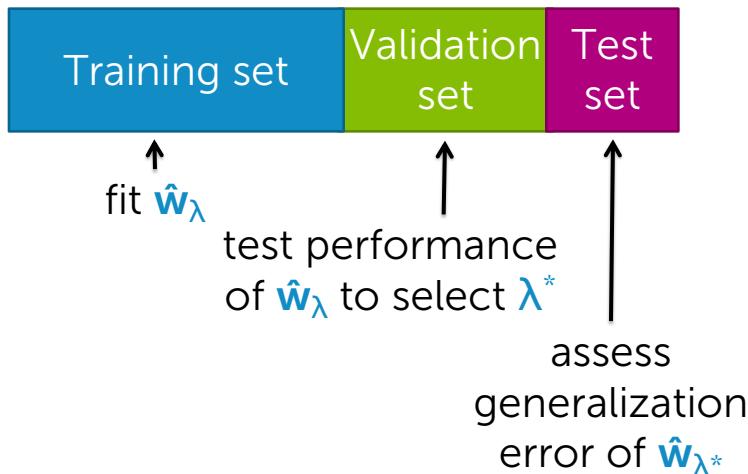
Ridge regression
(a.k.a L_2 regularization)

©2024 Emily Fox

CS 229: Machine Learning

6

How to choose tuning parameter

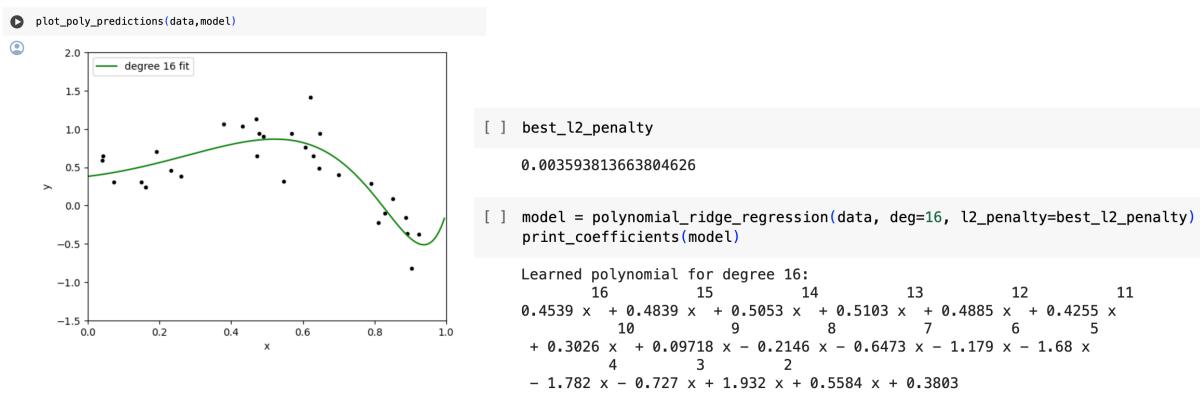


©2024 Emily Fox

CS 229: Machine Learning

7

Ridge penalty for polynomial regression (D=16)



©2024 Emily Fox

CS 229: Machine Learning

8

Lasso Regression:

Regularization for feature selection

CS 229: Machine Learning

Emily Fox

Stanford University

January 22, 2024

©2024 Emily Fox

9

Feature selection task

©2024 Emily Fox

CS 229: Machine Learning

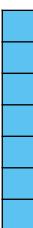
10

Why might you want to perform feature selection?

Efficiency:

- If size(w) = 100B, each prediction is expensive
- If \hat{w} [sparse], computation only depends on # of non-zeros

many zeros

$$\hat{y}_i = \sum_{\hat{w}_j \neq 0} \hat{w}_j h_j(x_i)$$


Interpretability:

- Which features are relevant for prediction?

©2024 Emily Fox

CS 229: Machine Learning

11

Sparsity: Housing application



| | |
|------------------------|------------------|
| Lot size | Dishwasher |
| Single Family | Garbage disposal |
| Year built | Microwave |
| Last sold price | Range / Oven |
| Last sale price/sqft | Refrigerator |
| Finished sqft | Washer |
| Unfinished sqft | Dryer |
| Finished basement sqft | Laundry location |
| # floors | Heating type |
| Flooring types | Jetted Tub |
| Parking type | Deck |
| Parking amount | Fenced Yard |
| Cooling | Lawn |
| Heating | Garden |
| Exterior materials | Sprinkler System |
| Roof type | : |
| Structure style | |

©2024 Emily Fox

CS 229: Machine Learning

12

Option 1: All subsets or greedy variants

©2024 Emily Fox

CS 229: Machine Learning

13

Exhaustive approach: “all subsets”

Consider all possible models, each using a subset of features

How many models were evaluated? each indexed by features included

| | | |
|--|-------------------------------------|--|
| $y_i = \epsilon_i$ | $[0 0 0 \dots 0 0 0]$ | $\left. \begin{matrix} \text{Select } 0 \\ \text{Select } 1 \\ \dots \\ \text{Select } D \end{matrix} \right\}$ |
| $y_i = w_0 h_0(\mathbf{x}_i) + \epsilon_i$ | $[1 0 0 \dots 0 0 0]$ | |
| $y_i = w_1 h_1(\mathbf{x}_i) + \epsilon_i$ | $[0 1 0 \dots 0 0 0]$ | |
| \vdots | \vdots | |
| $y_i = w_0 h_0(\mathbf{x}_i) + w_1 h_1(\mathbf{x}_i) + \dots + w_D h_D(\mathbf{x}_i) + \epsilon_i$ | $[1 1 0 \dots 0 0 0]$ | |
| \vdots | \vdots | |
| $y_i = w_0 h_0(\mathbf{x}_i) + w_1 h_1(\mathbf{x}_i) + \dots + w_D h_D(\mathbf{x}_i) + \epsilon_i$ | $[1 1 1 \dots 1 1 1]$ | 2^{D+1} |
| | $\vdots \vdots \vdots \dots \vdots$ | $2^8 = 256$ $2^{30} = 1,073,741,824$ $2^{1000} = 1.071509 \times 10^{301}$ $2^{100B} = \text{HUGE!!!!!!}$ |

Typically,
computationally
infeasible

©2024 Emily Fox

CS 229: Machine Learning

14

Choosing model complexity?

Option 1: Assess on validation set

Option 2: Cross validation *will see soon*

Option 3+: Other metrics for penalizing model complexity
like BIC...

©2024 Emily Fox

CS 229: Machine Learning

15

Greedy algorithms

Forward stepwise:

Starting from simple model and iteratively add features most useful to fit

Backward stepwise:

Start with full model and iteratively remove features least useful to fit

Combining forward and backward steps:

In forward algorithm, insert steps to remove features no longer as important

Lots of other variants, too.

©2024 Emily Fox

CS 229: Machine Learning

16

Option 2: Regularize

©2024 Emily Fox

CS 229: Machine Learning

17

Using regularization for feature selection

Instead of searching over a **discrete** set of solutions, can we use **regularization**?

- Start with full model (all possible features)
- “Shrink” some coefficients **exactly to 0**
 - i.e., knock out certain features
- Non-zero coefficients indicate “selected” features

©2024 Emily Fox

CS 229: Machine Learning

18

Ridge regression: L_2 regularized regression

Total cost =

$$\underbrace{\text{measure of fit}}_{\text{RSS}(\mathbf{w})} + \lambda \underbrace{\text{measure of magnitude of coefficients}}_{\|\mathbf{w}\|_2^2 = w_0^2 + \dots + w_D^2}$$

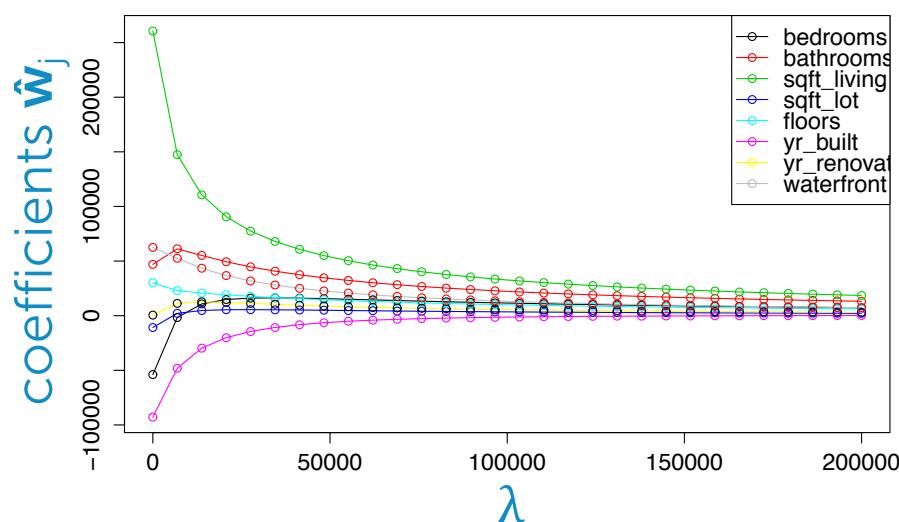
Encourages small weights
but not exactly 0

©2024 Emily Fox

CS 229: Machine Learning

19

Coefficient path – ridge



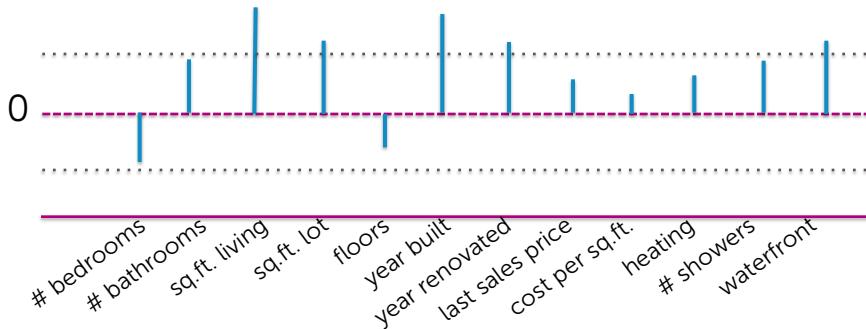
©2024 Emily Fox

CS 229: Machine Learning

20

Thresholding ridge coefficients?

Why don't we just set small ridge coefficients to 0?



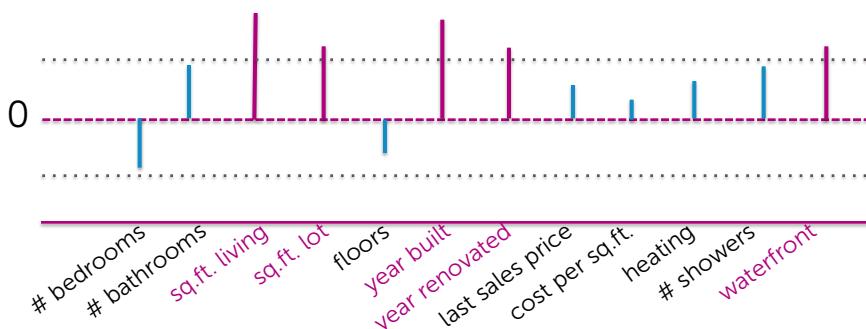
©2024 Emily Fox

CS 229: Machine Learning

21

Thresholding ridge coefficients?

Selected features for a given threshold value



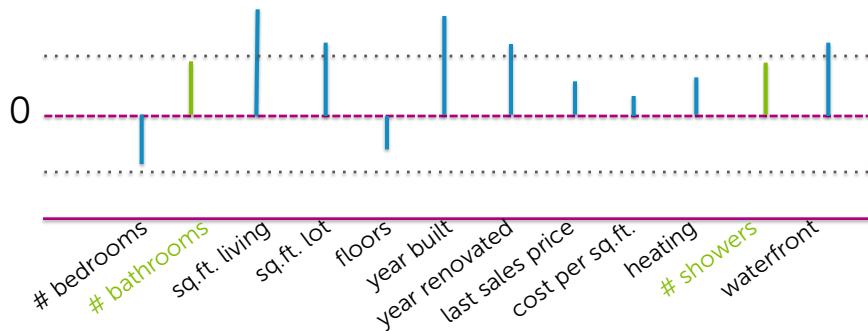
©2024 Emily Fox

CS 229: Machine Learning

22

Thresholding ridge coefficients?

Let's look at two related features...



Nothing measuring bathrooms was included!

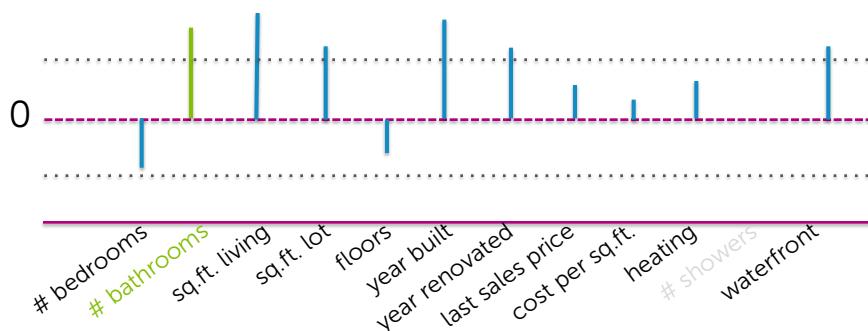
©2024 Emily Fox

CS 229: Machine Learning

23

Thresholding ridge coefficients?

If only one of the features had been included...



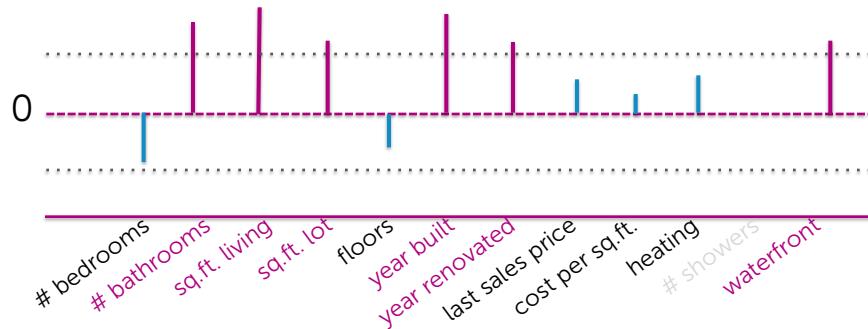
©2024 Emily Fox

CS 229: Machine Learning

24

Thresholding ridge coefficients?

Would have included bathrooms in selected model



Can regularization lead directly to sparsity?

©2024 Emily Fox

CS 229: Machine Learning

25

Try this cost instead of ridge...

Total cost =

$$\text{measure of fit} + \lambda \text{ measure of magnitude of coefficients}$$

RSS}(\mathbf{w})
||\mathbf{w}||_1 = |w_0| + \dots + |w_D|

Leads to sparse solutions!

Lasso regression
(a.k.a. L_1 regularized regression)

©2024 Emily Fox

CS 229: Machine Learning

26

Lasso regression: L_1 regularized regression

Just like ridge regression, solution is governed by a continuous parameter λ

$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$$

If $\lambda = 0$: $\hat{\mathbf{w}}^{\text{lasso}} = \hat{\mathbf{w}}^{\text{LS}}$ (unreg. soln.)

If $\lambda = \infty$: $\hat{\mathbf{w}}^{\text{lasso}} = \mathbf{0}$

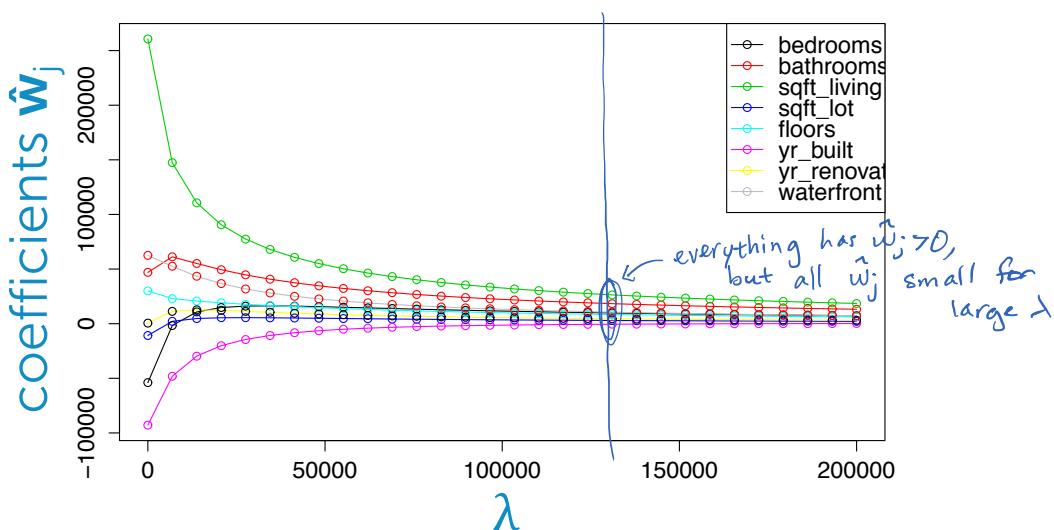
If λ in between: $0 \leq \|\hat{\mathbf{w}}^{\text{lasso}}\|_1 \leq \|\hat{\mathbf{w}}^{\text{LS}}\|_1$

©2024 Emily Fox

CS 229: Machine Learning

27

Coefficient path – ridge

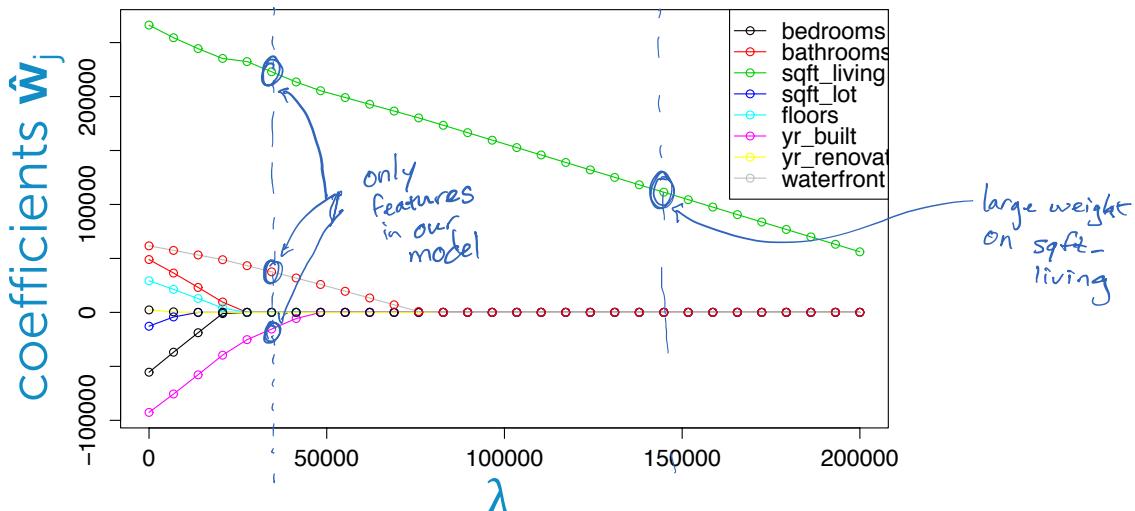


©2024 Emily Fox

CS 229: Machine Learning

28

Coefficient path – lasso



©2024 Emily Fox

CS 229: Machine Learning

29

Revisit polynomial fit demo

What happens if we refit our high-order polynomial, but now using **lasso regression**?

Will consider a few settings of λ ...

©2024 Emily Fox

CS 229: Machine Learning

30

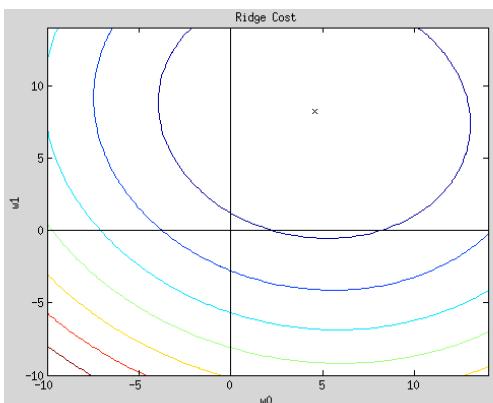
Geometric intuition for sparsity of lasso solution

©2024 Emily Fox

CS 229: Machine Learning

31

Visualizing the ridge cost in 2D



$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$$

↑ circles
weighting

Movie: fun of increasing λ
 \mathbf{x} = opt soln for a specific λ

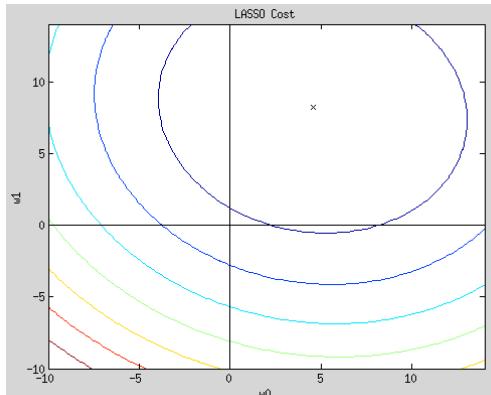
$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 = \sum_{i=1}^N (y_i - w_0 h_0(\mathbf{x}_i) - w_1 h_1(\mathbf{x}_i))^2 + \lambda (w_0^2 + w_1^2)$$

©2024 Emily Fox

CS 229: Machine Learning

32

Visualizing the lasso cost in 2D



$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$,
ellipses diamonds

$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1 = \sum_{i=1}^N (y_i - w_0 h_0(\mathbf{x}_i) - w_1 h_1(\mathbf{x}_i))^2 + \lambda (|w_0| + |w_1|)$$

©2024 Emily Fox

CS 229: Machine Learning

33

Fitting the lasso regression model
(for given λ value)

©2024 Emily Fox

CS 229: Machine Learning

34

How we optimized past objectives

To solve for $\hat{\mathbf{w}}$, previously took gradient of total cost objective and either:

- 1) Derived closed-form solution
- 2) Used in gradient descent algorithm

©2024 Emily Fox

CS 229: Machine Learning

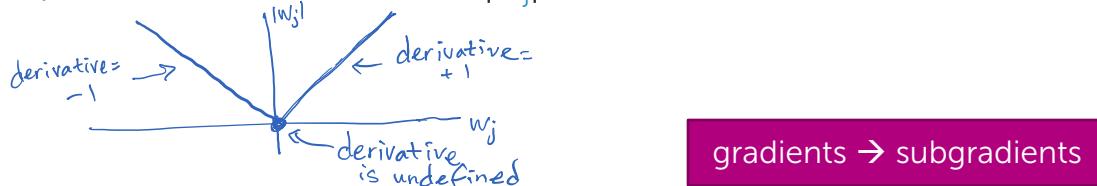
35

Optimizing the lasso objective

Lasso total cost: $\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$

Issues:

- 1) What's the derivative of $|w_j|$?



- 2) Even if we could compute derivative, no closed-form solution

can use subgradient descent

©2024 Emily Fox

CS 229: Machine Learning

36

Aside 1: Coordinate descent

©2024 Emily Fox

CS 229: Machine Learning

37

Coordinate descent

Goal: Minimize some function g

$$g(\mathbf{w}) = g(w_0, w_1, \dots, w_D)$$

$$\min_{\mathbf{w}} g(\mathbf{w})$$

when keeping others fixed

Often, hard to find minimum for all coordinates, but easy for each coordinate

Coordinate descent:

Initialize $\hat{\mathbf{w}} = 0$ (or smartly...)

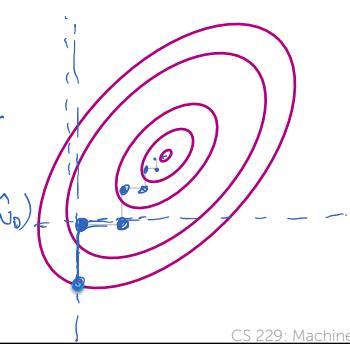
while not converged

pick a coordinate j

$$\hat{w}_j \leftarrow \min_{w_j} g(\hat{w}_0, \hat{w}_1, \dots, \hat{w}_{j-1}, w_j, \hat{w}_{j+1}, \dots, \hat{w}_D)$$

©2024 Emily Fox

values from prev. iter
just min j^{th} coord



CS 229: Machine Learning

38

19

Comments on coordinate descent

How do we pick next coordinate?

- At random ("random" or "stochastic" coordinate descent), round robin, ...

No stepsize to choose!

Super useful approach for *many* problems

- Converges to optimum in some cases
(e.g., "strongly convex")
- Converges for lasso objective

©2024 Emily Fox

CS 229: Machine Learning

39

Aside 2: Normalizing features

©2024 Emily Fox

CS 229: Machine Learning

40

20

Normalizing features

Scale training columns (**not rows!**) as:

$$\underline{h}_j(\mathbf{x}_k) = \frac{h_j(\mathbf{x}_k)}{\sqrt{\sum_{i=1}^N h_j(\mathbf{x}_i)^2}}$$

Normalizer:
 Z_j

Apply same training scale factors to test data:

$$\underline{h}_j(\mathbf{x}_k) = \frac{h_j(\mathbf{x}_k)}{\sqrt{\sum_{i=1}^N h_j(\mathbf{x}_i)^2}}$$

Normalizer:
 Z_j

*apply to
test point*

summing over training points



©2024 Emily Fox

CS 229: Machine Learning

41

Aside 3: Coordinate descent for unregularized regression (for normalized features)

©2024 Emily Fox

CS 229: Machine Learning

42

Optimizing least squares objective one coordinate at a time

$$\text{RSS}(\mathbf{w}) = \sum_{i=1}^N \left(y_i - \sum_{j=0}^D w_j h_j(\mathbf{x}_i) \right)^2$$

normalized feature

Fix all coordinates \mathbf{w}_{-j} and take partial w.r.t. w_j

$$\begin{aligned} \frac{\partial}{\partial w_j} \text{RSS}(\mathbf{w}) &= -2 \sum_{i=1}^N h_j(\mathbf{x}_i) \left(y_i - \sum_{k=0}^D w_k h_k(\mathbf{x}_i) \right) \\ &= -2 \sum_{i=1}^N h_j \left(y_i - \sum_{k \neq j} w_k h_k \right) + 2 w_j \sum_{i=1}^N h_j^2 \\ &\triangleq p_j \\ &= -2 p_j + 2 w_j \end{aligned}$$

by defn of normalized feat.
= 1

©2024 Emily Fox

CS 229: Machine Learning

43

Optimizing least squares objective one coordinate at a time

$$\text{RSS}(\mathbf{w}) = \sum_{i=1}^N \left(y_i - \sum_{j=0}^D w_j h_j(\mathbf{x}_i) \right)^2$$

Set partial = 0 and solve

$$\begin{aligned} \frac{\partial}{\partial w_j} \text{RSS}(\mathbf{w}) &= -2 p_j + 2 \hat{w}_j = 0 \\ \hat{w}_j &= p_j \end{aligned}$$

©2024 Emily Fox

CS 229: Machine Learning

44

Coordinate descent for least squares regression

Initialize $\hat{\mathbf{w}} = 0$ (or smartly...)

while not converged

for $j=0,1,\dots,D$

$$\text{compute: } \rho_j = \sum_{i=1}^N h_j(\mathbf{x}_i)(y_i - \hat{y}_i(\hat{\mathbf{w}}_{-j}))$$

↑ residual
without feature j
↑ prediction without feature j

set: $\hat{\mathbf{w}}_j = \rho_j$

©2024 Emily Fox

CS 229: Machine Learning

45

Coordinate descent for lasso (for normalized features)

©2024 Emily Fox

CS 229: Machine Learning

46

Coordinate descent for least squares regression

Initialize $\hat{\mathbf{w}} = 0$ (or smartly...)

while not converged

for $j=0,1,\dots,D$

$$\text{compute: } \rho_j = \sum_{i=1}^N h_j(\mathbf{x}_i)(y_i - \hat{y}_i(\hat{\mathbf{w}}_{-j}))$$

↑
residual
without feature j

↑
set: $\hat{w}_j = \rho_j$ prediction without feature j

©2024 Emily Fox

CS 229: Machine Learning

47

Coordinate descent for lasso

Initialize $\hat{\mathbf{w}} = 0$ (or smartly...)

while not converged

for $j=0,1,\dots,D$

$$\text{compute: } \rho_j = \sum_{i=1}^N h_j(\mathbf{x}_i)(y_i - \hat{y}_i(\hat{\mathbf{w}}_{-j}))$$

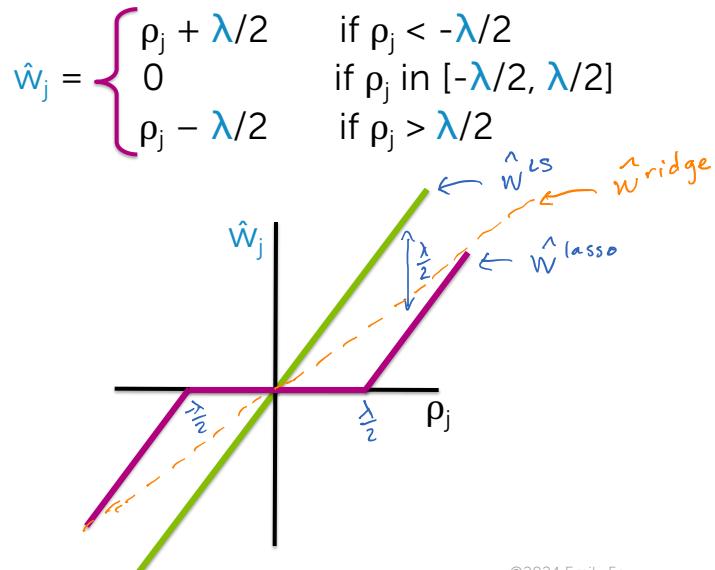
↑
set: $\hat{w}_j = \begin{cases} \rho_j + \lambda/2 & \text{if } \rho_j < -\lambda/2 \\ 0 & \text{if } \rho_j \in [-\lambda/2, \lambda/2] \\ \rho_j - \lambda/2 & \text{if } \rho_j > \lambda/2 \end{cases}$

©2024 Emily Fox

CS 229: Machine Learning

48

Soft thresholding



©2024 Emily Fox

CS 229: Machine Learning

49

How to assess convergence?

Initialize $\hat{\mathbf{w}} = 0$ (or smartly...)

while not converged

for $j=0,1,\dots,D$

compute: $\rho_j = \sum_{i=1}^N h_j(\mathbf{x}_i)(y_i - \hat{y}_i(\hat{\mathbf{w}}_{-j}))$

set: $\hat{w}_j = \begin{cases} \rho_j + \lambda/2 & \text{if } \rho_j < -\lambda/2 \\ 0 & \text{if } \rho_j \in [-\lambda/2, \lambda/2] \\ \rho_j - \lambda/2 & \text{if } \rho_j > \lambda/2 \end{cases}$

©2024 Emily Fox

CS 229: Machine Learning

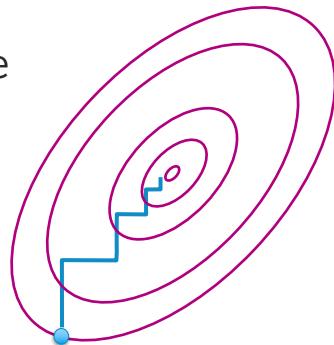
50

Convergence criteria

When to stop?

For convex problems, will start to take
smaller and smaller steps

Measure size of steps taken in a
full loop over all features
– stop when $\text{max step} < \epsilon$



©2024 Emily Fox

CS 229: Machine Learning

51

Other lasso solvers

Classically: Least angle regression ([LARS](#)) [[Efron et al. '04](#)]

Then: [Coordinate descent](#) algorithm [[Fu '98](#), [Friedman, Hastie, & Tibshirani '08](#)]

Now:

- [Parallel CD](#) (e.g., [Shotgun](#), [[Bradley et al. '11](#)])
- Other parallel learning approaches for linear models
 - Parallel stochastic gradient descent ([SGD](#)) (e.g., [Hogwild!](#) [[Niu et al. '11](#)])
 - Parallel independent solutions then [averaging](#) [[Zhang et al. '12](#)]
- Alternating directions method of multipliers ([ADMM](#)) [[Boyd et al. '11](#)]

©2024 Emily Fox

CS 229: Machine Learning

52

Coordinate descent for lasso (for unnormalized features)

©2024 Emily Fox

CS 229: Machine Learning

53

Coordinate descent for lasso with unnormalized features

Precompute: $z_j = \sum_{i=1}^N h_j(\mathbf{x}_i)^2$

Initialize $\hat{\mathbf{w}} = 0$ (or smartly...)

while not converged

for $j=0,1,\dots,D$

compute: $\rho_j = \sum_{i=1}^N h_j(\mathbf{x}_i)(y_i - \hat{y}_i(\hat{\mathbf{w}}_{-j}))$

set: $\hat{w}_j = \begin{cases} (\rho_j + \lambda/2)/z_j & \text{if } \rho_j < -\lambda/2 \\ 0 & \text{if } \rho_j \in [-\lambda/2, \lambda/2] \\ (\rho_j - \lambda/2)/z_j & \text{if } \rho_j > \lambda/2 \end{cases}$

©2024 Emily Fox

CS 229: Machine Learning

54

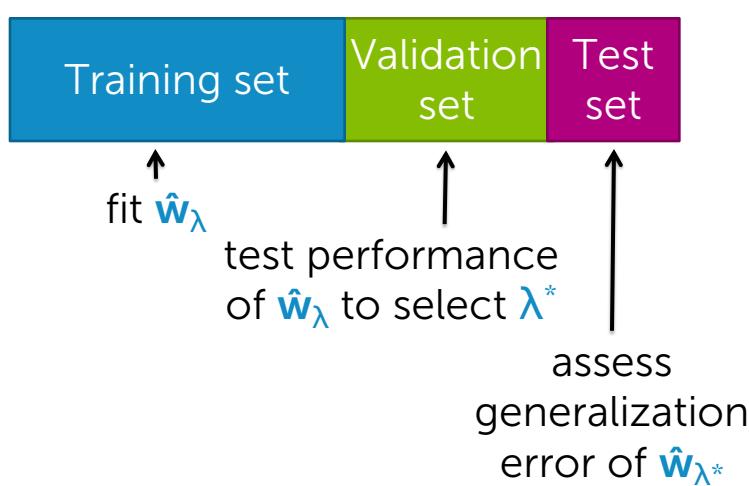
How to choose λ

©2024 Emily Fox

CS 229: Machine Learning

55

If sufficient amount of data...



©2024 Emily Fox

CS 229: Machine Learning

56

Start with smallish dataset

All data

©2024 Emily Fox

CS 229: Machine Learning

57

Still form test set and hold out

Rest of data

Test
set

©2024 Emily Fox

CS 229: Machine Learning

58

29

How do we use the other data?

Rest of data

use for both training and validation, but not so naively

©2024 Emily Fox

CS 229: Machine Learning

59

Recall naïve approach

Training set

Valid.
set

↑
small validation set

Is validation set enough to compare performance of \hat{w}_λ across λ values?

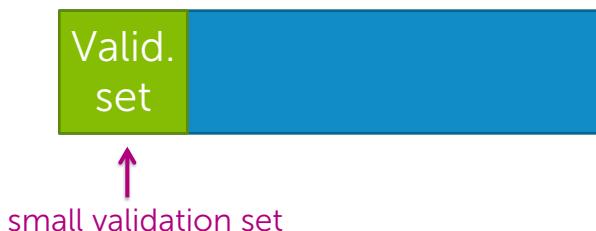
No

©2024 Emily Fox

CS 229: Machine Learning

60

Choosing the validation set



Didn't have to use the last data points tabulated to form validation set

Can use **any** data subset

©2024 Emily Fox

CS 229: Machine Learning

61

Choosing the validation set



Which subset should I use?

ALL!

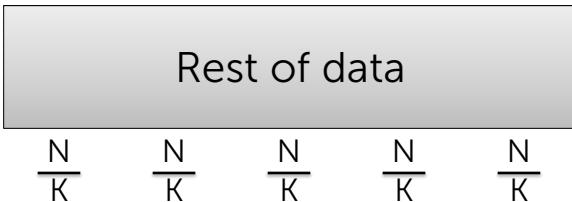
average performance
over all choices

©2024 Emily Fox

CS 229: Machine Learning

62

K-fold cross validation



Preprocessing: Randomly assign data to K groups

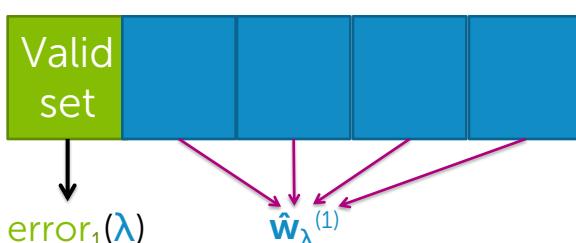
(use same split of data for all other steps)

©2024 Emily Fox

CS 229: Machine Learning

63

K-fold cross validation



For $k=1, \dots, K$

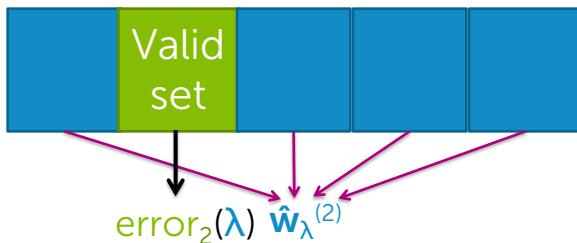
1. Estimate $\hat{w}_\lambda^{(k)}$ on the training blocks
2. Compute error on validation block: $\text{error}_k(\lambda)$

©2024 Emily Fox

CS 229: Machine Learning

64

K-fold cross validation



For $k=1, \dots, K$

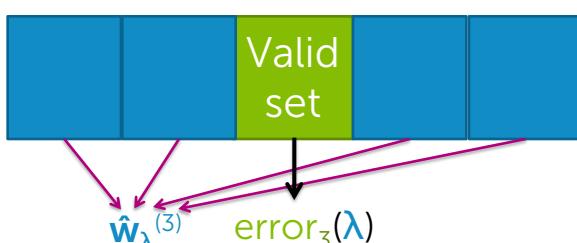
1. Estimate $\hat{w}_\lambda^{(k)}$ on the training blocks
2. Compute error on validation block: $\text{error}_k(\lambda)$

©2024 Emily Fox

CS 229: Machine Learning

65

K-fold cross validation



For $k=1, \dots, K$

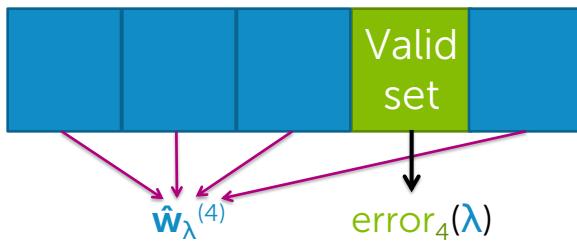
1. Estimate $\hat{w}_\lambda^{(k)}$ on the training blocks
2. Compute error on validation block: $\text{error}_k(\lambda)$

©2024 Emily Fox

CS 229: Machine Learning

66

K-fold cross validation



For $k=1, \dots, K$

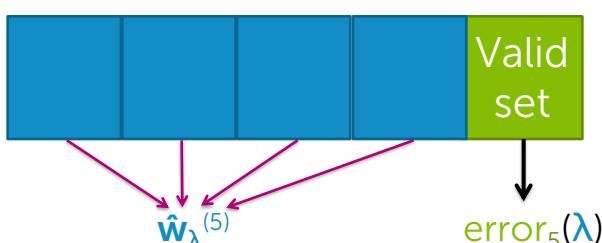
1. Estimate $\hat{w}_{\lambda}^{(k)}$ on the training blocks
2. Compute error on validation block: $\text{error}_k(\lambda)$

©2024 Emily Fox

CS 229: Machine Learning

67

K-fold cross validation



For $k=1, \dots, K$

1. Estimate $\hat{w}_{\lambda}^{(k)}$ on the training blocks
2. Compute error on validation block: $\text{error}_k(\lambda)$

Compute average error: $\text{CV}(\lambda) = \frac{1}{K} \sum_{k=1}^K \text{error}_k(\lambda)$

©2024 Emily Fox

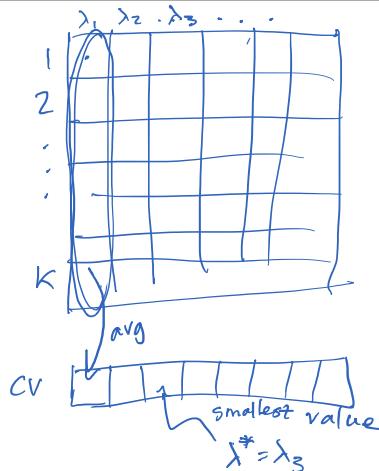
CS 229: Machine Learning

68

K-fold cross validation



Repeat procedure for each choice of λ



Choose λ^* to minimize $CV(\lambda)$

©2024 Emily Fox

CS 229: Machine Learning

69

What value of K?

Formally, the **best approximation** occurs for validation sets of size 1 ($K=N$)

leave-one-out
cross validation

Computationally intensive

- requires computing N fits of model per λ

Typically, $K=5$ or 10

5-fold CV

10-fold CV

©2024 Emily Fox

CS 229: Machine Learning

70

Choosing λ via cross validation for lasso

Cross validation is choosing the λ that provides best predictive accuracy

Tends to favor less sparse solutions, and thus smaller λ , than optimal choice for feature selection

extensions developed that are selection consistent

c.f., "Probabilistic Machine Learning: An Introduction", Murphy, 2022 for further discussion

©2024 Emily Fox

CS 229: Machine Learning

71

Practical concerns with lasso

PRACTICALITIES

©2024 Emily Fox

CS 229: Machine Learning

72

Debiasing lasso

Lasso shrinks coefficients relative to LS solution
 → more bias, less variance

Can reduce bias as follows:

1. Run lasso to select features
2. Run least squares regression with only selected features

"Relevant" features no longer shrunk relative to LS fit of same reduced model

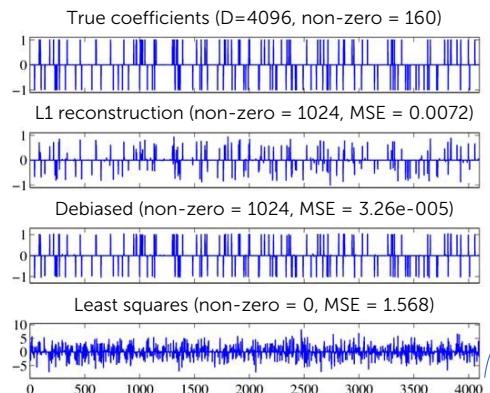


Figure used with permission of Mario Figueiredo
 (captions modified to fit course)

©2024 Emily Fox

CS 229: Machine Learning

73

Issues with standard lasso objective

1. With group of highly correlated features, lasso tends to select amongst them arbitrarily
 - Often prefer to select all together
2. Often, empirically ridge has better predictive performance than lasso, but lasso leads to sparser solution

Elastic net aims to address these issues

- hybrid between lasso and ridge regression
- uses L_1 and L_2 penalties

See **Zou & Hastie '05** for further discussion

©2024 Emily Fox

CS 229: Machine Learning

74

Summary for feature selection and lasso regression

©2024 Emily Fox

CS 229: Machine Learning

75

What you can do now...

- Describe “all subsets” and greedy variants for feature selection
- Analyze computational costs of these algorithms
- Formulate lasso objective
- Describe what happens to estimated lasso coefficients as tuning parameter λ is varied
- Interpret lasso coefficient path plot
- Contrast ridge and lasso regression
- Estimate lasso regression parameters using an iterative coordinate descent algorithm
- Implement K-fold cross validation to select lasso tuning parameter λ

©2024 Emily Fox

CS 229: Machine Learning

76

38

Deriving the lasso coordinate descent update

OPTIONAL



©2024 Emily Fox

CS 229: Machine Learning

77

Optimizing lasso objective one coordinate at a time

$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1 = \sum_{i=1}^N \left(y_i - \sum_{j=0}^D w_j h_j(\mathbf{x}_i) \right)^2 + \lambda \sum_{j=0}^D |w_j|$$

Fix all coordinates \mathbf{w}_{-j} and take partial w.r.t. w_j

derive *without* normalizing features

©2024 Emily Fox

CS 229: Machine Learning

78

Part 1: Partial of RSS term

$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1 = \sum_{i=1}^N (y_i - \sum_{j=0}^D w_j h_j(\mathbf{x}_i))^2 + \lambda \sum_{j=0}^D |w_j|$$

$$\begin{aligned} \frac{\partial}{\partial w_j} \text{RSS}(\mathbf{w}) &= -2 \sum_{i=1}^N h_j(\mathbf{x}_i) (y_i - \sum_{k=0}^D w_k h_k(\mathbf{x}_i)) \\ &= -2 \sum_{i=1}^N h_j(\mathbf{x}_i) \left(y_i - \sum_{k \neq j} w_k h_k(\mathbf{x}_i) - w_j h_j(\mathbf{x}_i) \right) \\ &= -2 \sum_{i=1}^N h_j(\mathbf{x}_i) (y_i - \underbrace{\sum_{k \neq j} w_k h_k(\mathbf{x}_i)}_{\stackrel{\triangle}{=} p_j}) + 2 w_j \sum_{i=1}^N h_j(\mathbf{x}_i)^2 \\ &= -2 p_j + 2 w_j z_j \end{aligned}$$

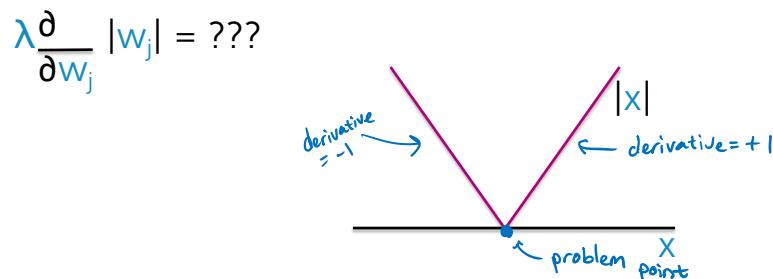
©2024 Emily Fox

CS 229: Machine Learning

79

Part 2: Partial of L_1 penalty term

$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1 = \sum_{i=1}^N (y_i - \sum_{j=0}^D w_j h_j(\mathbf{x}_i))^2 + \lambda \sum_{j=0}^D |w_j|$$



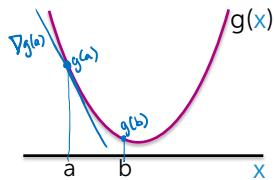
©2024 Emily Fox

CS 229: Machine Learning

80

Subgradients of convex functions

Gradients lower bound convex functions:

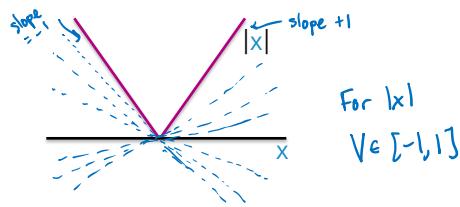


$$g(b) \geq g(a) + \underline{\nabla g(a)(b-a)}$$

unique at x if function differentiable at x

Subgradients: Generalize gradients to non-differentiable points:

- Any plane that lower bounds function



For $|x|$
 $\forall x \in [-1, 1]$

$\nabla g(x)$ subgradient of g at x
 if
 $g(b) \geq g(a) + \nabla g(a)(b-a)$

©2024 Emily Fox

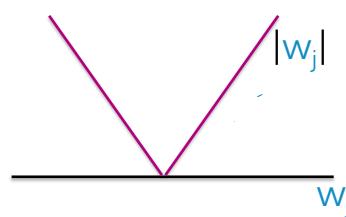
CS 229: Machine Learning

81

Part 2: Subgradient of L_1 term

$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1 = \sum_{i=1}^N (y_i - \sum_{j=0}^D w_j h_j(\mathbf{x}_i))^2 + \lambda \sum_{j=0}^D |w_j|$$

$$\lambda \partial_{w_j} |w_j| = \begin{cases} -\lambda & \text{when } w_j < 0 \\ [-\lambda, \lambda] & \text{when } w_j = 0 \\ \lambda & \text{when } w_j > 0 \end{cases}$$



©2024 Emily Fox

CS 229: Machine Learning

82

Putting it all together...

$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1 = \sum_{i=1}^N (y_i - \sum_{j=0}^D w_j h_j(\mathbf{x}_i))^2 + \lambda \sum_{j=0}^D |w_j|$$

$$\partial_{w_j} [\text{lasso cost}] = 2z_j w_j - 2\rho_j + \begin{cases} -\lambda & \text{when } w_j < 0 \\ [-\lambda, \lambda] & \text{when } w_j = 0 \\ \lambda & \text{when } w_j > 0 \end{cases}$$

$$= \begin{cases} 2z_j w_j - 2\rho_j - \lambda & \text{when } w_j < 0 \\ [-2\rho_j - \lambda, -2\rho_j + \lambda] & \text{when } w_j = 0 \\ 2z_j w_j - 2\rho_j + \lambda & \text{when } w_j > 0 \end{cases}$$

©2024 Emily Fox

CS 229: Machine Learning

83

Optimal solution:
Set subgradient = 0

$$\partial_{w_j} [\text{lasso cost}] = \begin{cases} 2z_j w_j - 2\rho_j - \lambda & \text{when } w_j < 0 \\ [-2\rho_j - \lambda, -2\rho_j + \lambda] & \text{when } w_j = 0 \\ 2z_j w_j - 2\rho_j + \lambda & \text{when } w_j > 0 \end{cases} = 0$$

Case 1 ($w_j < 0$): $2z_j \hat{w}_j - 2\rho_j - \lambda = 0$ For $\hat{w}_j < 0$, need $\rho_j < -\frac{\lambda}{2}$

$$\hat{w}_j = \frac{2\rho_j + \lambda}{2z_j} = \frac{\rho_j + \frac{\lambda}{2}}{z_j}$$

Case 2 ($w_j = 0$): $\hat{w}_j = 0$ For $\hat{w}_j = 0$, need $[-2\rho_j - \lambda, -2\rho_j + \lambda]$ to contain 0:
 $-2\rho_j + \lambda \geq 0 \rightarrow \rho_j \leq \frac{\lambda}{2}$ so that $\hat{w}_j = 0$ is an optimum
 $-2\rho_j - \lambda \leq 0 \rightarrow \rho_j \geq -\frac{\lambda}{2}$ $\frac{\lambda}{2} \leq \rho_j \leq -\frac{\lambda}{2}$

Case 3 ($w_j > 0$): $2z_j \hat{w}_j - 2\rho_j + \lambda = 0$ For $\hat{w}_j > 0$, need $\rho_j > \frac{\lambda}{2}$

$$\hat{w}_j = \frac{\rho_j - \frac{\lambda}{2}}{z_j}$$

©2024 Emily Fox

CS 229: Machine Learning

84

Optimal solution:
Set subgradient = 0

$$\partial_{w_j} [\text{lasso cost}] = 0 \quad \begin{cases} 2z_j w_j - 2\rho_j - \lambda & \text{when } w_j < 0 \\ [-2\rho_j - \lambda, -2\rho_j + \lambda] & \text{when } w_j = 0 \\ 2z_j w_j - 2\rho_j + \lambda & \text{when } w_j > 0 \end{cases}$$



$$\hat{w}_j = \begin{cases} (\rho_j + \lambda/2)/z_j & \text{if } \rho_j < -\lambda/2 \\ 0 & \text{if } \rho_j \in [-\lambda/2, \lambda/2] \\ (\rho_j - \lambda/2)/z_j & \text{if } \rho_j > \lambda/2 \end{cases}$$

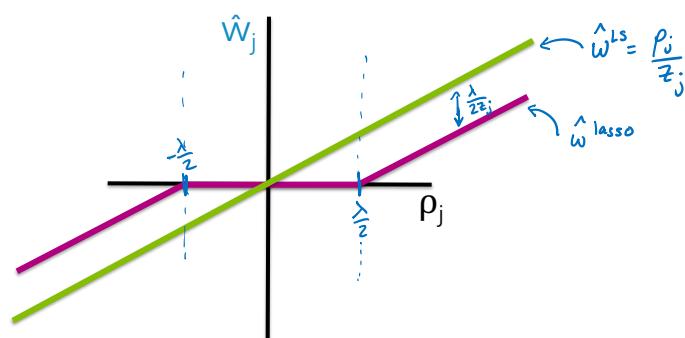
©2024 Emily Fox

CS 229: Machine Learning

85

Soft thresholding

$$\hat{w}_j = \begin{cases} (\rho_j + \lambda/2)/z_j & \text{if } \rho_j < -\lambda/2 \\ 0 & \text{if } \rho_j \in [-\lambda/2, \lambda/2] \\ (\rho_j - \lambda/2)/z_j & \text{if } \rho_j > \lambda/2 \end{cases}$$



©2024 Emily Fox

CS 229: Machine Learning

86

Coordinate descent for lasso

Precompute: $z_j = \sum_{i=1}^N h_j(\mathbf{x}_i)^2$

Initialize $\hat{\mathbf{w}} = 0$ (or smartly...)

while not converged

for $j=0,1,\dots,D$

compute: $\rho_j = \sum_{i=1}^N h_j(\mathbf{x}_i)(y_i - \hat{y}_i(\hat{\mathbf{w}}_{-j}))$

set: $\hat{w}_j = \begin{cases} (\rho_j + \lambda/2)/z_j & \text{if } \rho_j < -\lambda/2 \\ 0 & \text{if } \rho_j \text{ in } [-\lambda/2, \lambda/2] \\ (\rho_j - \lambda/2)/z_j & \text{if } \rho_j > \lambda/2 \end{cases}$

©2024 Emily Fox

CS 229: Machine Learning

87

Geometric intuition for sparsity of lasso solution in more detail

OPTIONAL

©2024 Emily Fox

CS 229: Machine Learning

88

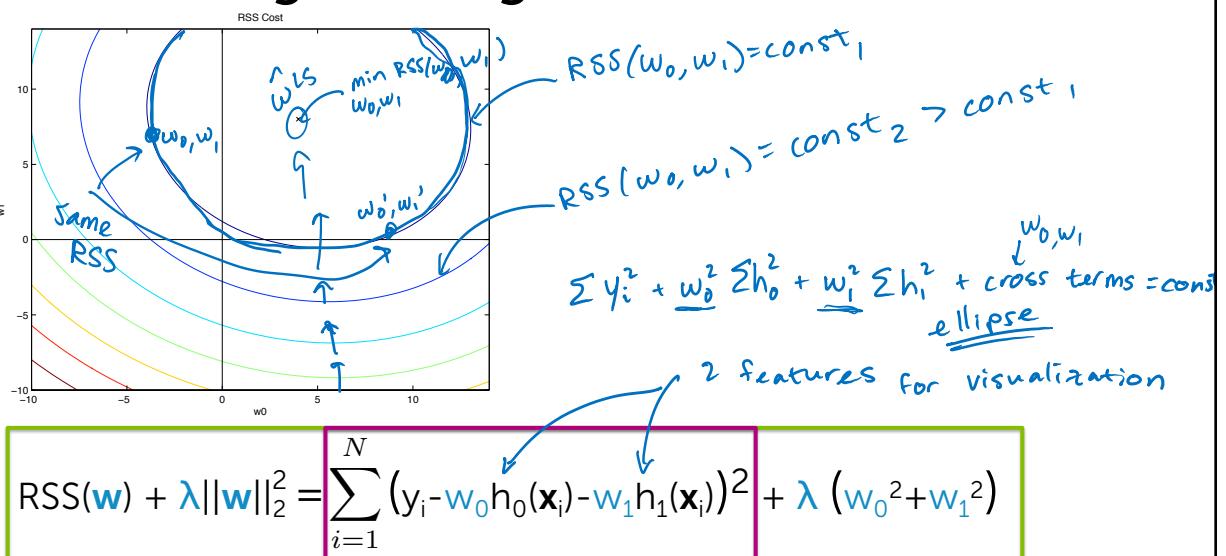
Geometric intuition for ridge regression

©2024 Emily Fox

CS 229: Machine Learning

89

Visualizing the ridge cost in 2D

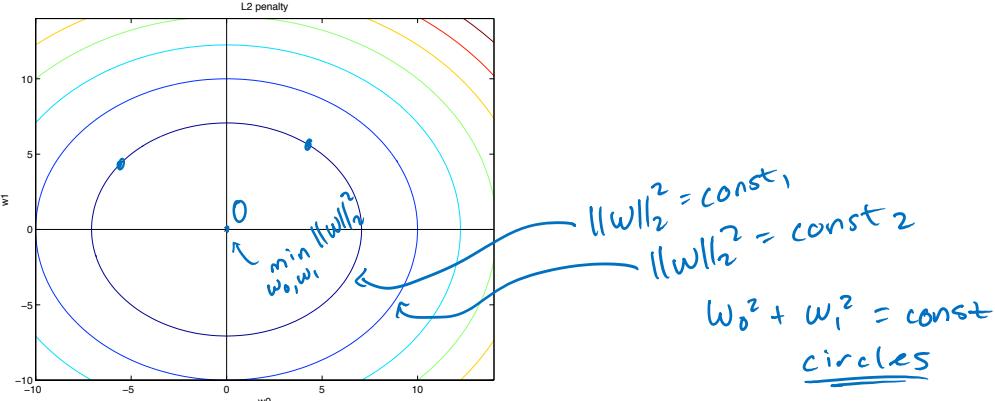


©2024 Emily Fox

CS 229: Machine Learning

90

Visualizing the ridge cost in 2D



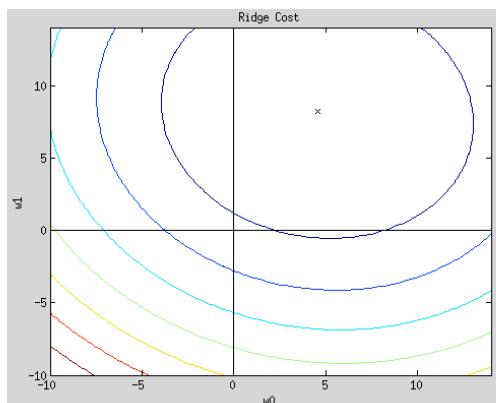
$$\text{RSS}(\mathbf{w}) + \lambda \| \mathbf{w} \|_2^2 = \sum_{i=1}^N (y_i - w_0 h_0(\mathbf{x}_i) - w_1 h_1(\mathbf{x}_i))^2 + \lambda (w_0^2 + w_1^2)$$

©2024 Emily Fox

CS 229: Machine Learning

91

Visualizing the ridge cost in 2D



Add contour plots together
 $\text{RSS}(\mathbf{w}) + \lambda \| \mathbf{w} \|_2^2$
 ellipses \uparrow circles
 weighting

Movie: fcn of increasing λ
 $x = \text{opt soln for specific } \lambda$

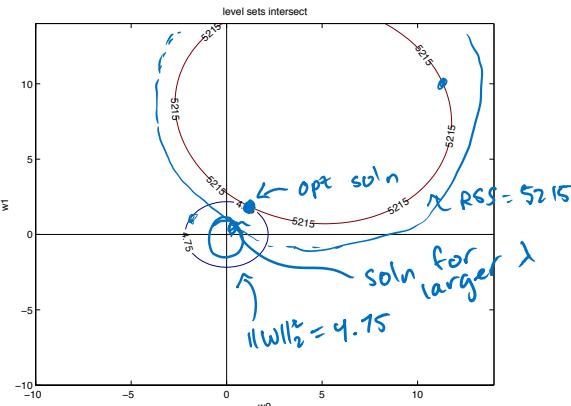
$$\text{RSS}(\mathbf{w}) + \lambda \| \mathbf{w} \|_2^2 = \sum_{i=1}^N (y_i - w_0 h_0(\mathbf{x}_i) - w_1 h_1(\mathbf{x}_i))^2 + \lambda (w_0^2 + w_1^2)$$

©2024 Emily Fox

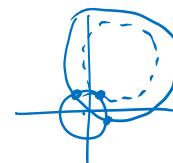
CS 229: Machine Learning

92

Visualizing the ridge solution



For specific λ value,
some balance between
RSS and $\|w\|_2^2$



$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 = \sum_{i=1}^N (y_i - \mathbf{w}_0 h_0(\mathbf{x}_i) - \mathbf{w}_1 h_1(\mathbf{x}_i))^2 + \lambda (\mathbf{w}_0^2 + \mathbf{w}_1^2)$$

©2024 Emily Fox

CS 229: Machine Learning

93

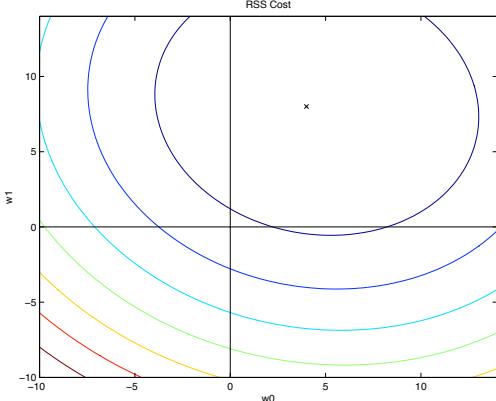
Geometric intuition for lasso

©2024 Emily Fox

CS 229: Machine Learning

94

Visualizing the lasso cost in 2D



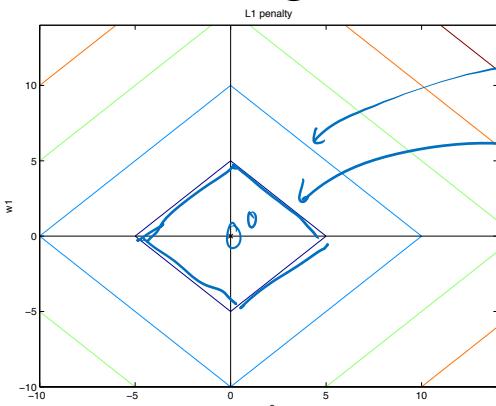
$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1 = \sum_{i=1}^N (y_i - w_0 h_0(\mathbf{x}_i) - w_1 h_1(\mathbf{x}_i))^2 + \lambda (|w_0| + |w_1|)$$

©2024 Emily Fox

CS 229: Machine Learning

95

Visualizing the lasso cost in 2D



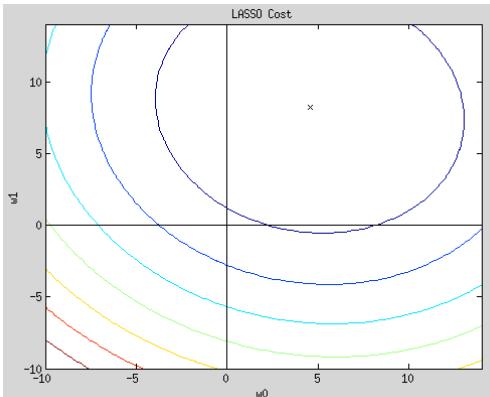
$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1 = \sum_{i=1}^N (y_i - w_0 h_0(\mathbf{x}_i) - w_1 h_1(\mathbf{x}_i))^2 + \lambda (|w_0| + |w_1|)$$

©2024 Emily Fox

CS 229: Machine Learning

96

Visualizing the lasso cost in 2D



$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$,
ellipses diamonds

Movie: as before

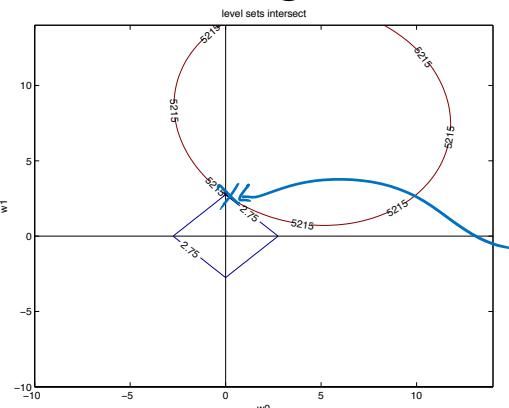
$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1 = \sum_{i=1}^N (y_i - w_0 h_0(x_i) - w_1 h_1(x_i))^2 + \lambda (|w_0| + |w_1|)$$

©2024 Emily Fox

CS 229: Machine Learning

97

Visualizing the lasso solution



solution for
specific value
of λ

$$\text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1 = \sum_{i=1}^N (y_i - w_0 h_0(x_i) - w_1 h_1(x_i))^2 + \lambda (|w_0| + |w_1|)$$

©2024 Emily Fox

CS 229: Machine Learning

98