

Lasso Regression:

Regularization for feature selection

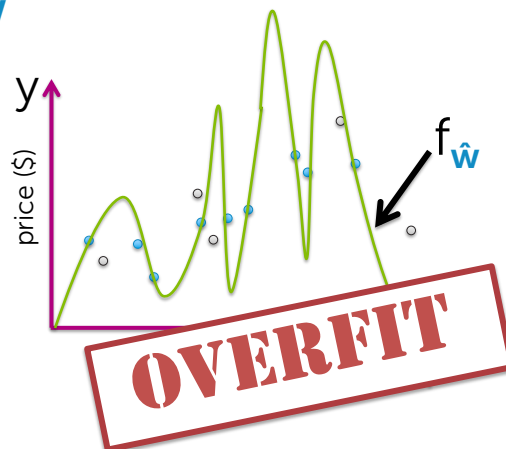
CSE 416: Machine Learning
Emily Fox
University of Washington
April 12, 2018

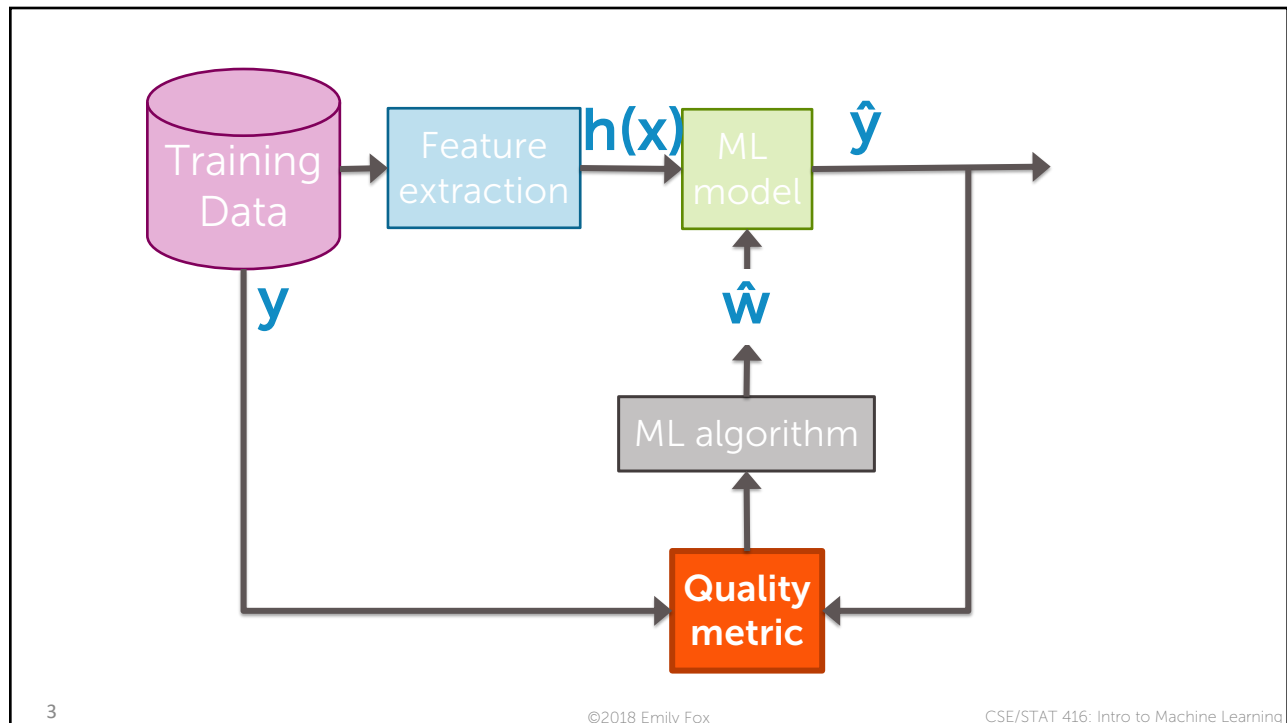
©2018 Emily Fox

Symptom of overfitting

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

Very large
coefficients (+/-)





Consider specific total cost

Want to balance:

- How well function fits data
- Magnitude of coefficients

Total cost =

$$\underbrace{\text{measure of fit}}_{\text{RSS}(\mathbf{w})} + \underbrace{\text{measure of magnitude of coefficients}}_{\|\mathbf{w}\|_2^2 = \sum_{j=0}^D w_j^2}$$

Consider resulting objective

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

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

λ tuning parameter = balance of fit and magnitude

Ridge regression
(a.k.a L_2 regularization)

5

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Measure of magnitude of regression coefficient

What summary # is indicative of size of regression coefficients?

- Sum? $w_0 = 1,527,301$ $w_1 = -1,605,253$
 $w_0 + w_1 = \text{small \#}$ X
 - Sum of absolute value? $\sum_{j=0}^D |w_j| \triangleq \|\mathbf{w}\|_1$
 - Sum of squares (L_2 norm) $\sum_{j=0}^D w_j^2 \triangleq \|\mathbf{w}\|_2^2$
- $\mathbf{w} = [w_0 \ w_1 \ \dots \ w_D]$
- L_1 norm... this discuss next lecture
- ridge \leftarrow focus of this lecture

6

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Feature selection task

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Why might you want to perform feature selection?

Efficiency:

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

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

pred. value → \hat{y}_i *est. coeff.* → \hat{w}_j *many zeros* → **sparse**

only sum over features w/ non-zero weights → $\hat{w}_j \neq 0$

Interpretability:

- Which features are relevant for prediction?

8

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Sparsity: Housing application



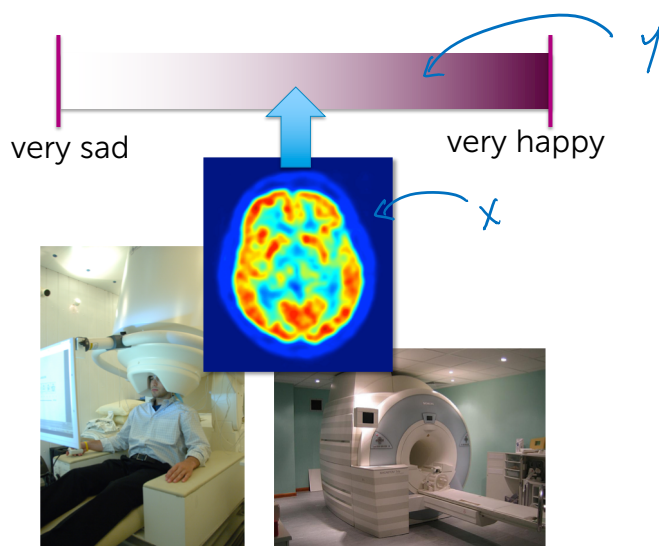
| | |
|------------------------|------------------|
| Lot size | Dishwasher |
| Single Family | Garbage disposal |
| Year built | Microwave |
| Last sold price | Range / Oven |
| Last sale price/sqft | Refrigerator |
| Finished sqft | Washer |
| Unfinished sqft | Dryer |
| Finished basement sqft | Laundry 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 | : |

9

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Sparsity: Reading your mind



Activity in which brain regions can predict happiness?

10

©2018 Emily Fox

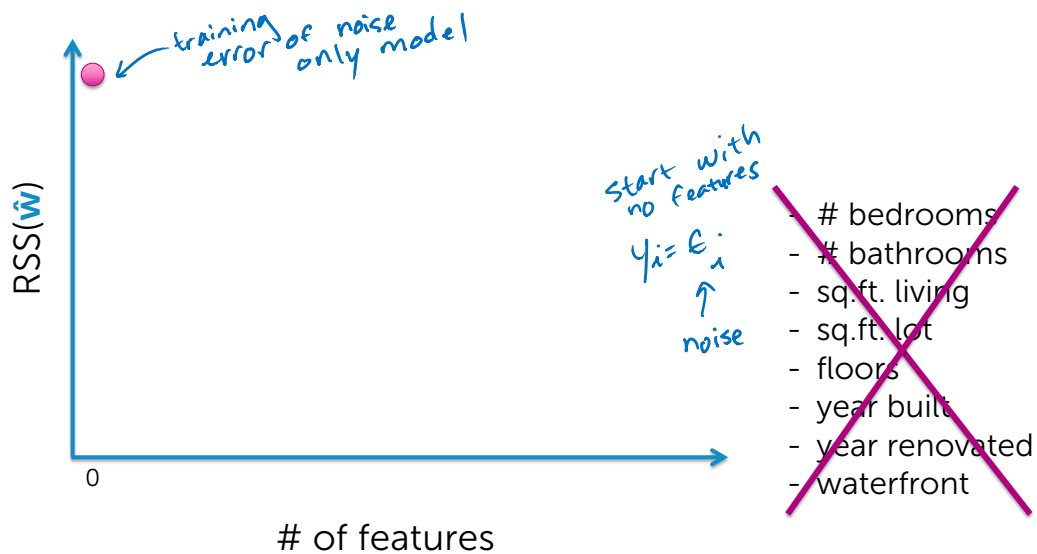
CSE/STAT 416: Intro to Machine Learning

Option 1: All subsets or greedy variants

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Find best model of size: 0

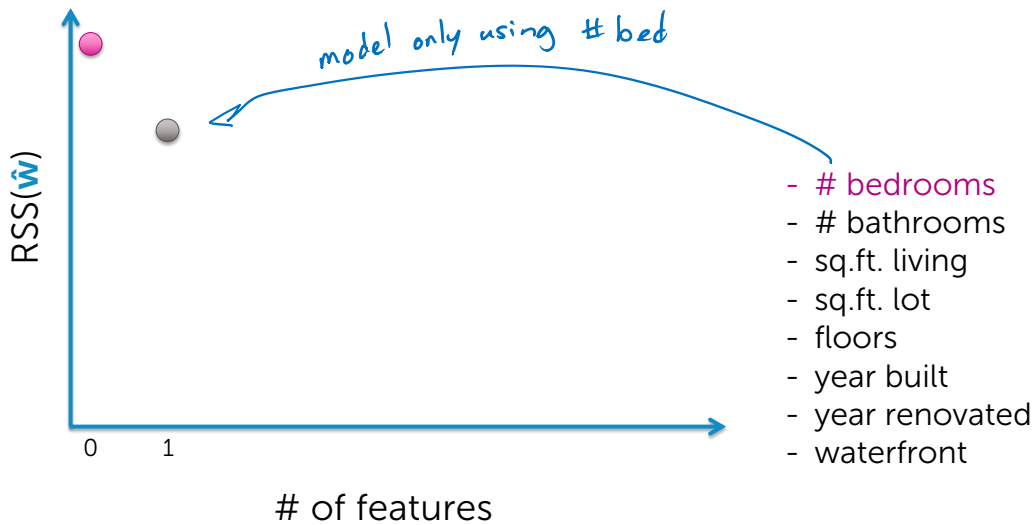


12

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Find best model of size: 1

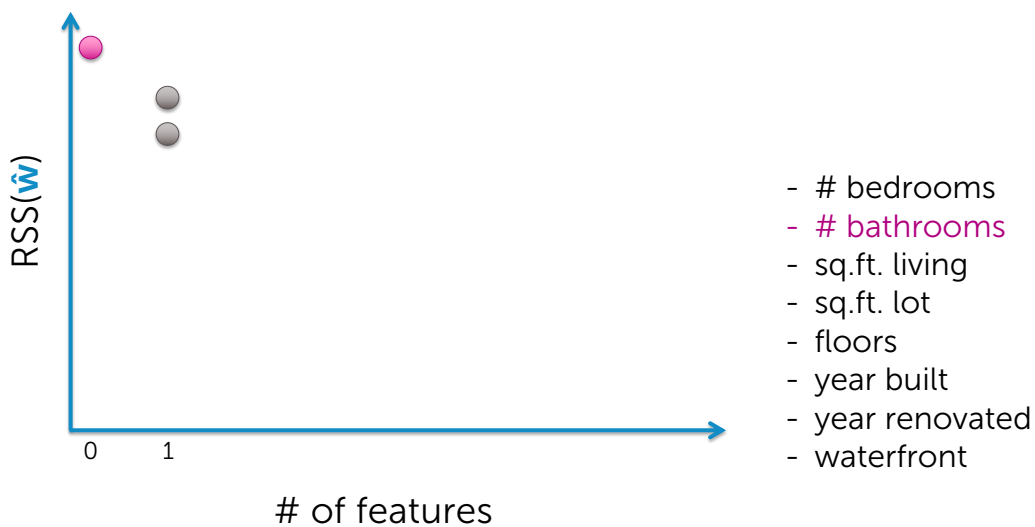


13

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Find best model of size: 1

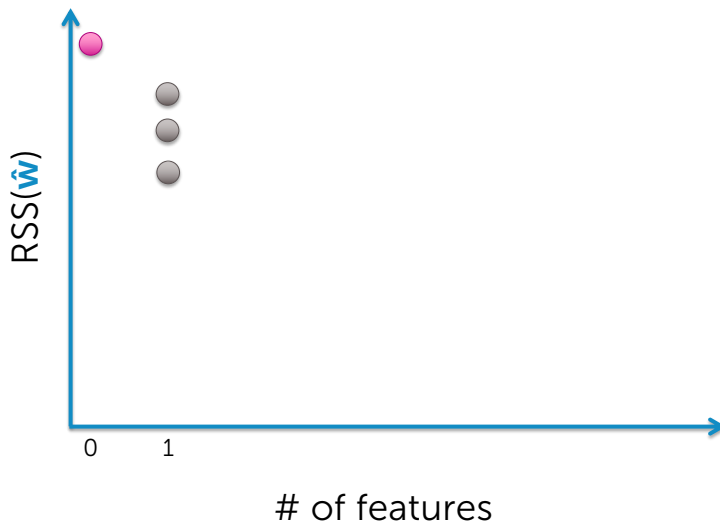


14

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Find best model of size: 1



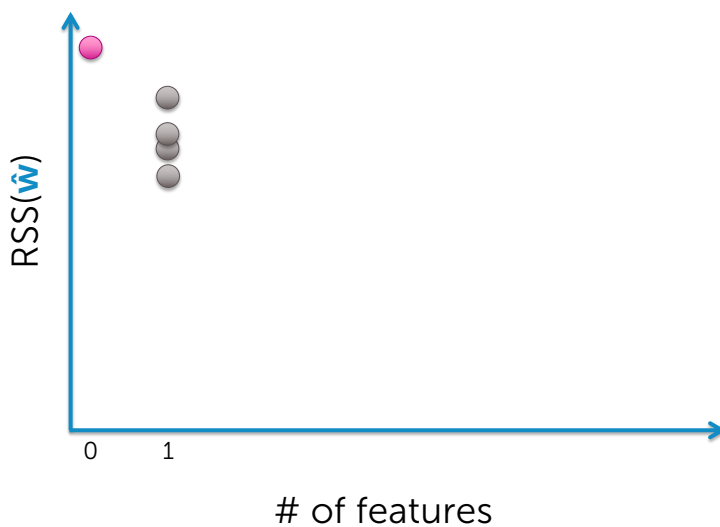
- # bedrooms
- # bathrooms
- sq.ft. living
- sq.ft. lot
- floors
- year built
- year renovated
- waterfront

15

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Find best model of size: 1



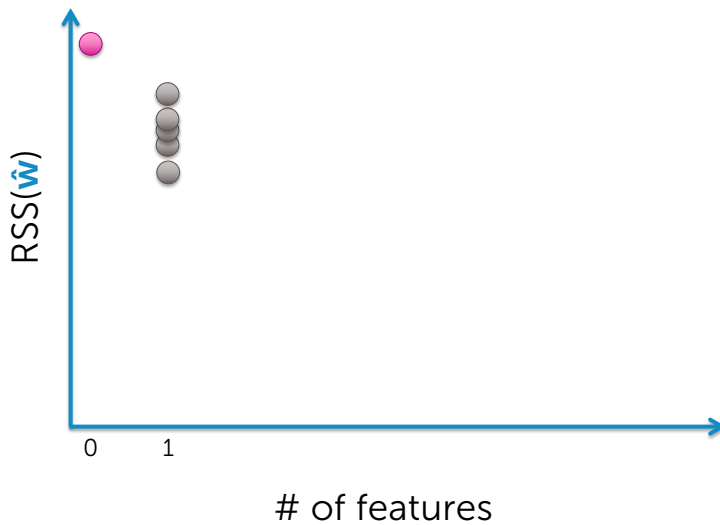
- # bedrooms
- # bathrooms
- sq.ft. living
- sq.ft. lot
- floors
- year built
- year renovated
- waterfront

16

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Find best model of size: 1



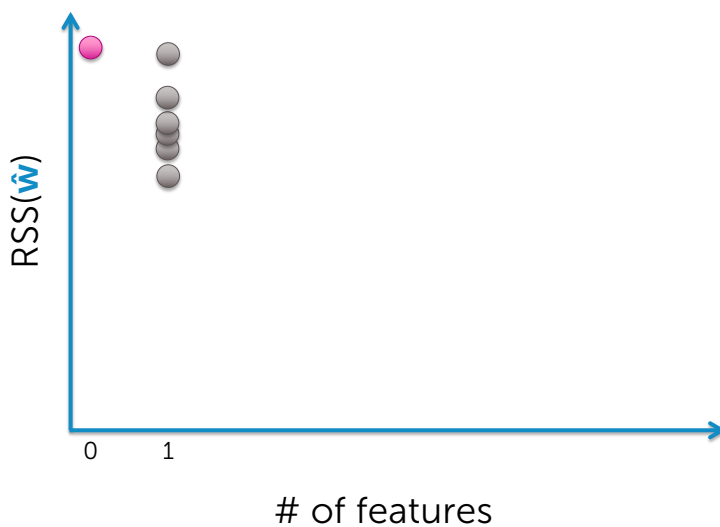
- # bedrooms
- # bathrooms
- sq.ft. living
- sq.ft. lot
- floors
- year built
- year renovated
- waterfront

17

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Find best model of size: 1



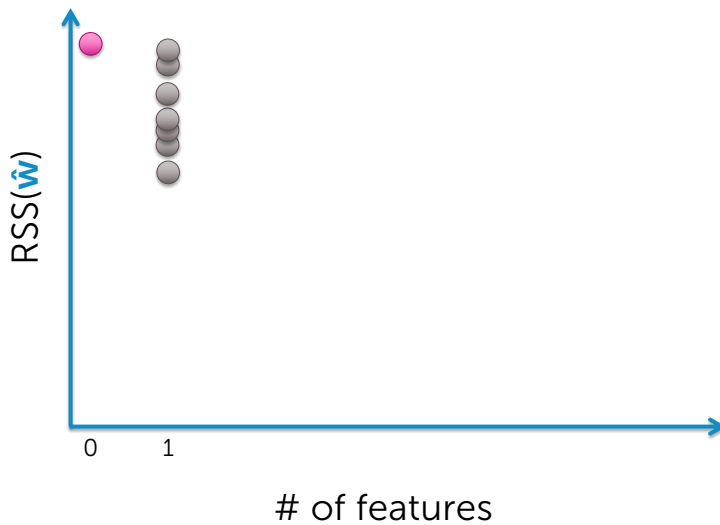
- # bedrooms
- # bathrooms
- sq.ft. living
- sq.ft. lot
- floors
- year built
- year renovated
- waterfront

18

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Find best model of size: 1



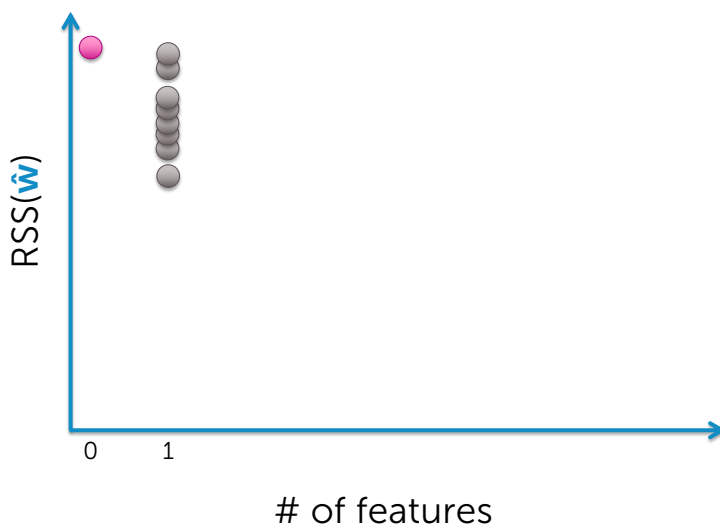
- # bedrooms
- # bathrooms
- sq.ft. living
- sq.ft. lot
- floors
- year built
- year renovated
- waterfront

19

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Find best model of size: 1



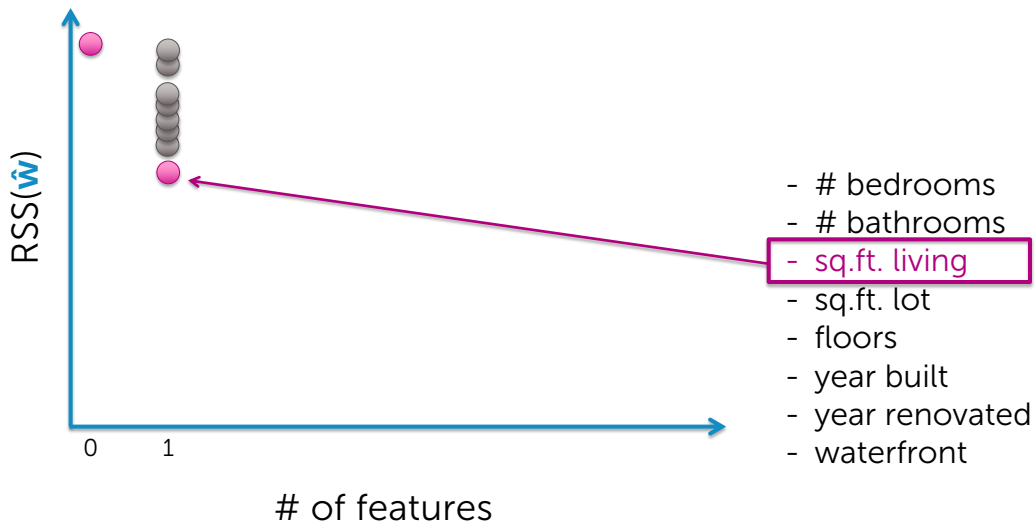
- # bedrooms
- # bathrooms
- sq.ft. living
- sq.ft. lot
- floors
- year built
- year renovated
- waterfront

20

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Find best model of size: 1



21

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Find best model of size: 2

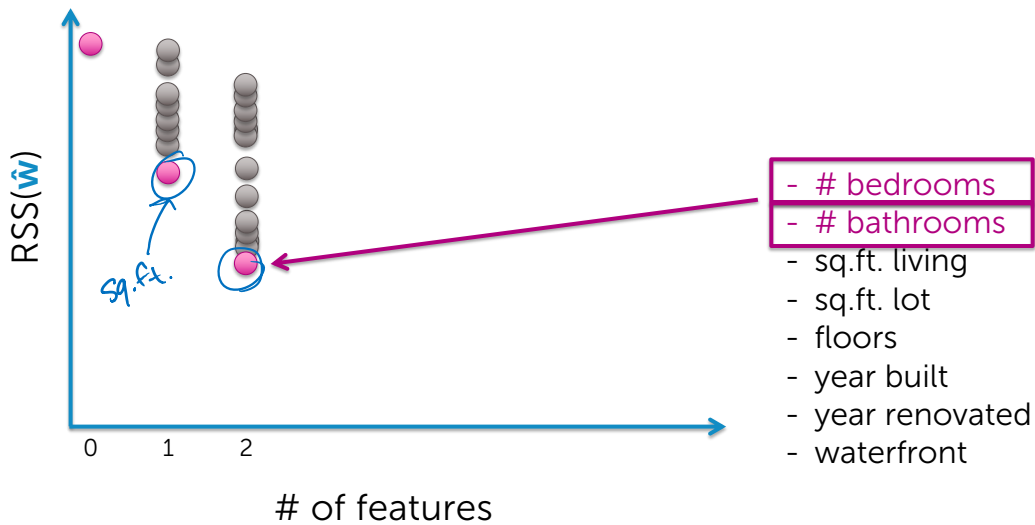


22

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Note: Not necessarily nested!

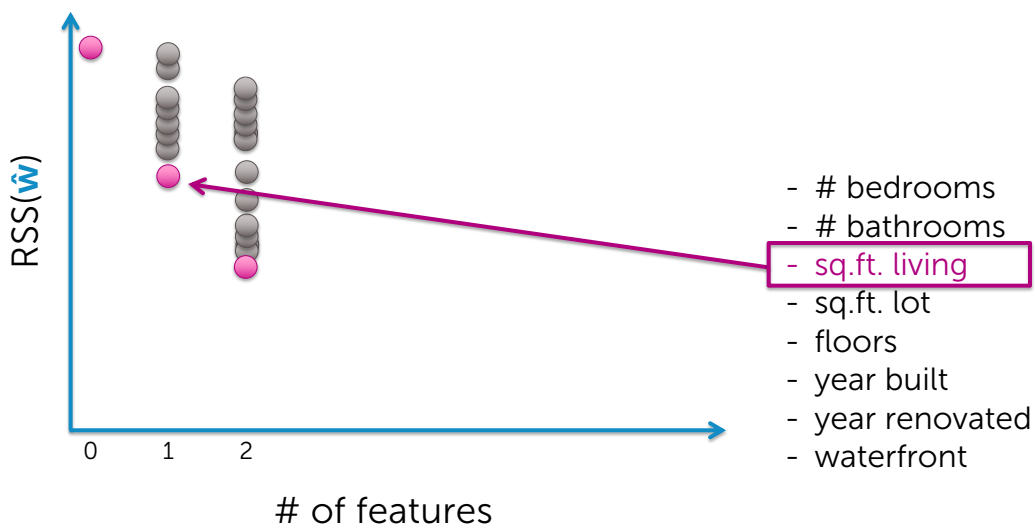


23

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Note: Not necessarily nested!

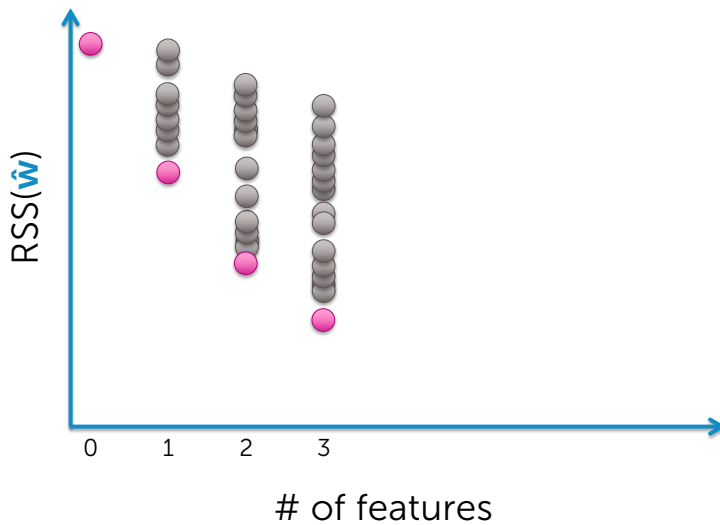


24

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Find best model of size: 3



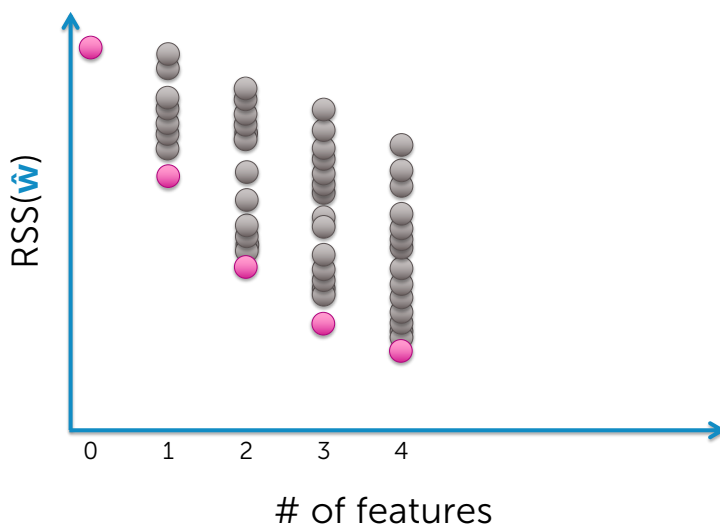
- # bedrooms
- # bathrooms
- sq.ft. living
- sq.ft. lot
- floors
- year built
- year renovated
- waterfront

25

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Find best model of size: 4



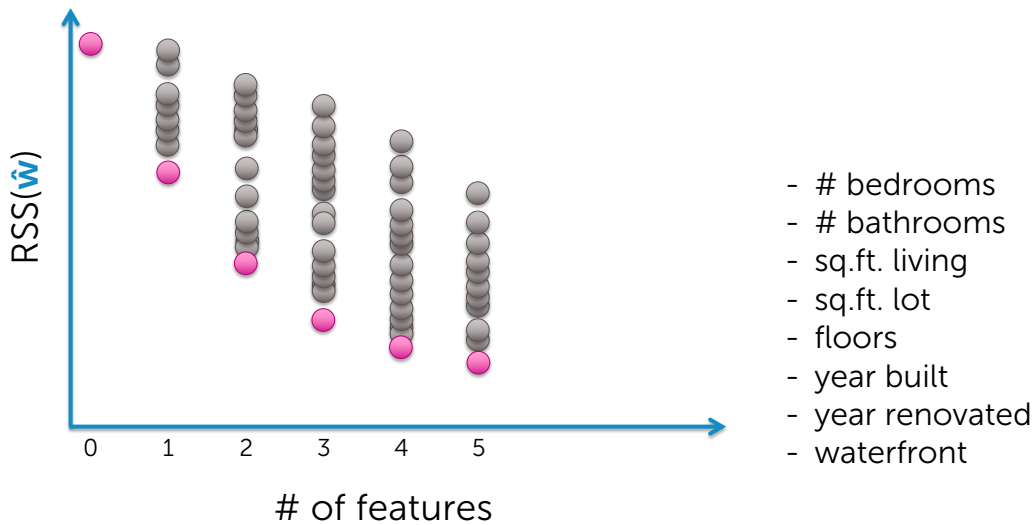
- # bedrooms
- # bathrooms
- sq.ft. living
- sq.ft. lot
- floors
- year built
- year renovated
- waterfront

26

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Find best model of size: 5

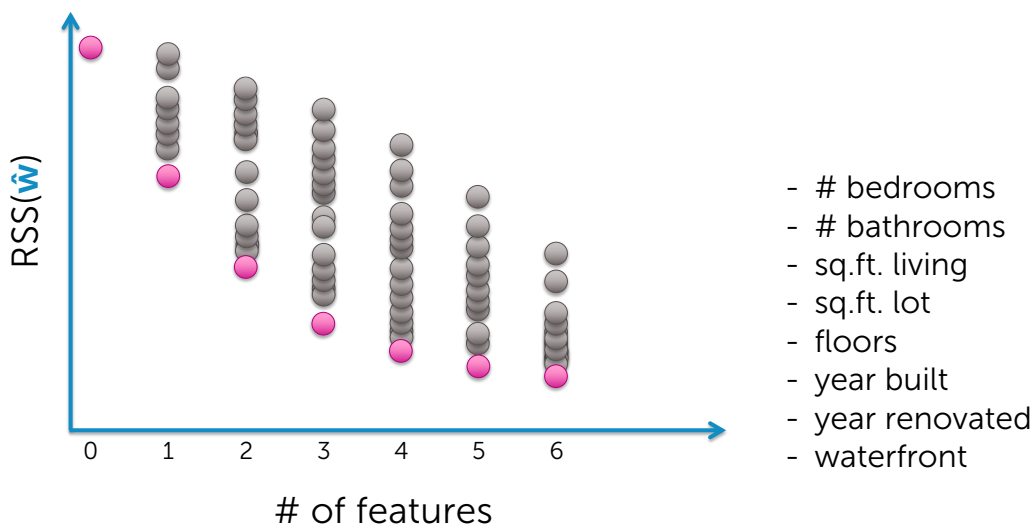


27

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Find best model of size: 6

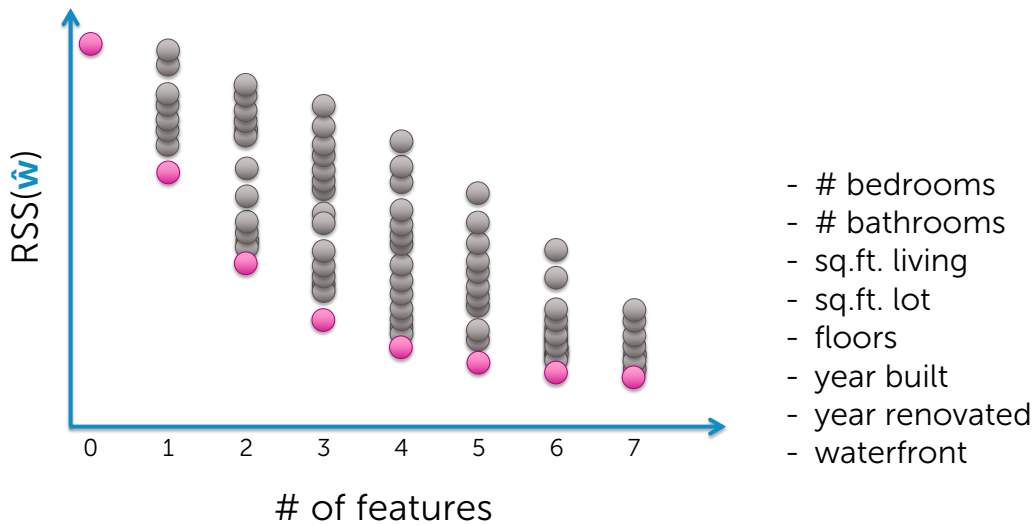


28

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Find best model of size: 7

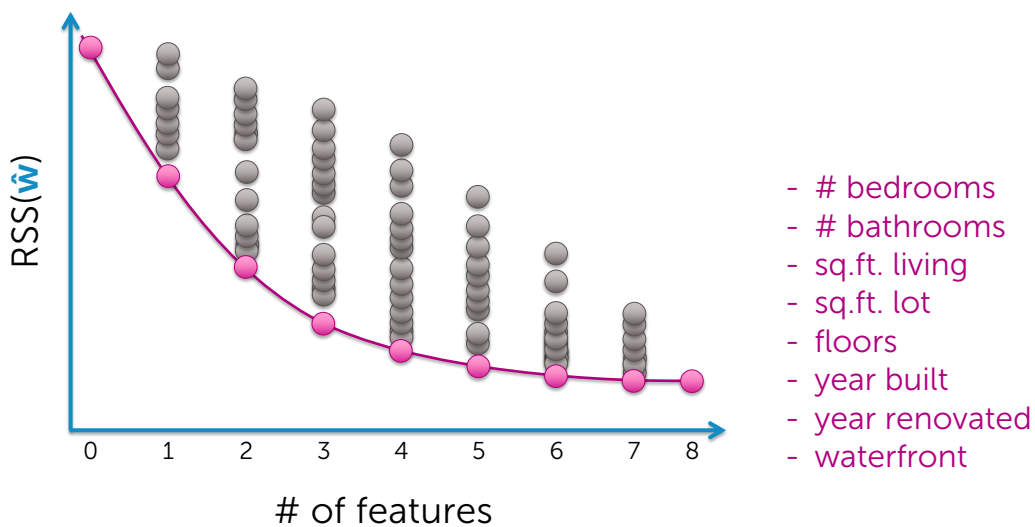


29

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Find best model of size: 8



30

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Choosing model complexity?

Option 1: Assess on validation set

Option 2: Cross validation

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

31

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Complexity of “all subsets”

How many models were evaluated?

– each indexed by features included

$$y_i = \varepsilon_i$$

$$y_i = w_0 h_0(\mathbf{x}_i) + \varepsilon_i$$

$$y_i = w_1 h_1(\mathbf{x}_i) + \varepsilon_i$$

$$\vdots$$

$$y_i = w_0 h_0(\mathbf{x}_i) + w_1 h_1(\mathbf{x}_i) + \varepsilon_i$$

$$\vdots$$

$$y_i = w_0 h_0(\mathbf{x}_i) + w_1 h_1(\mathbf{x}_i) + \dots + w_D h_D(\mathbf{x}_i) + \varepsilon_i$$

Feat 0 ← 0 if “no”
1 if “yes”
Feat 1 ... Feat D

0 0 0 ... 0 0 0
1 0 0 ... 0 0 0
0 1 0 ... 0 0 0
⋮
1 1 0 ... 0 0 0
⋮
1 1 1 ... 1 1 1
2 2 2 ... 2

$$\begin{aligned} 2^8 &= 256 \\ 2^{30} &= 1,073,741,824 \\ 2^{1000} &= 1.071509 \times 10^{301} \\ 2^{1008} &= \text{HUGE!!!!!!} \end{aligned}$$

 2^{D+1}

Typically,
computationally
infeasible

32

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

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.

33

©2017 Emily Fox

CSE/STAT 416: Intro to Machine Learning 33

Option 2: Regularize

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Ridge regression: L_2 regularized regression

Total cost =

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

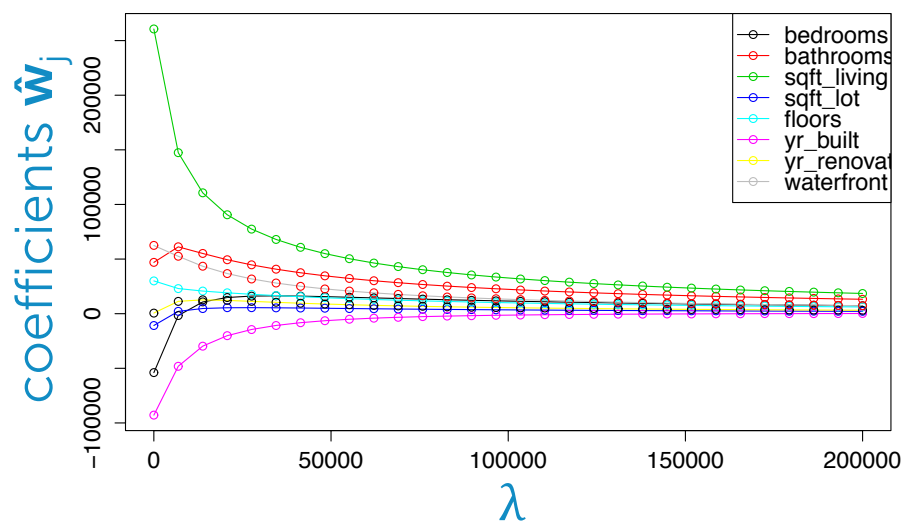
Encourages small weights
but not exactly 0

35

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Coefficient path – ridge



36

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

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

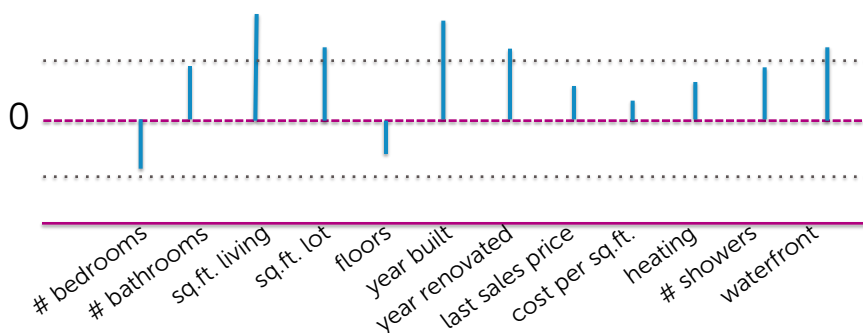
37

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Thresholding ridge coefficients?

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



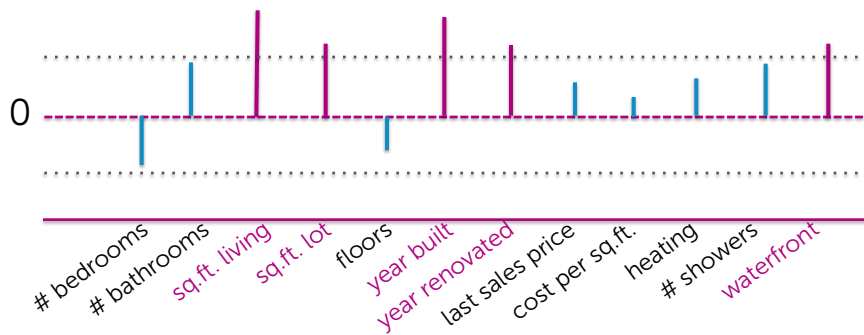
38

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Thresholding ridge coefficients?

Selected features for a given threshold value



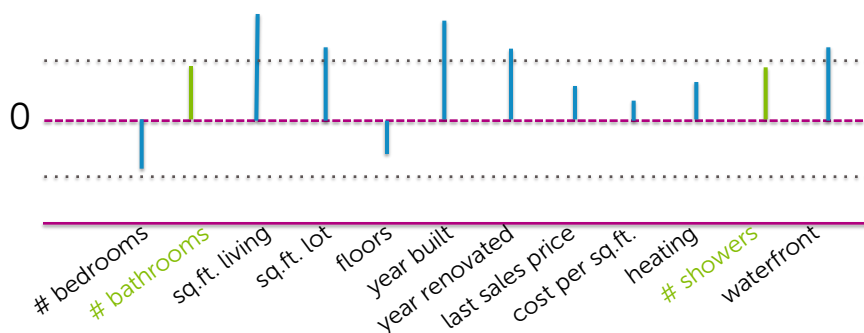
39

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Thresholding ridge coefficients?

Let's look at two related features...



Nothing measuring bathrooms was included!

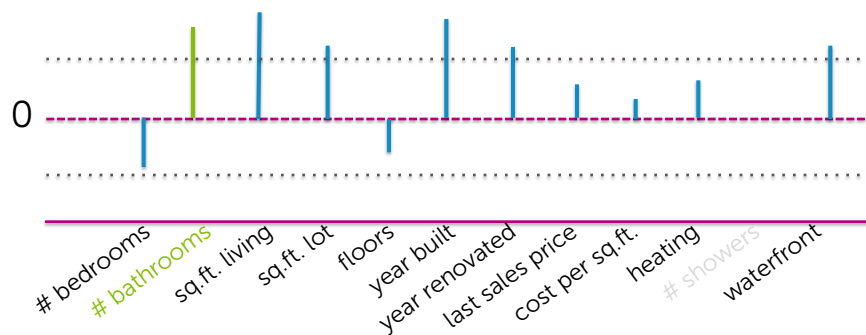
40

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Thresholding ridge coefficients?

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



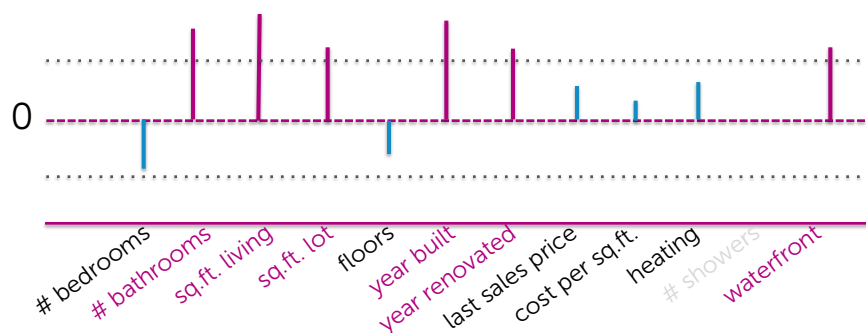
41

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Thresholding ridge coefficients?

Would have included bathrooms in selected model



Can regularization lead directly to sparsity?

42

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Try this cost instead of ridge...

Total cost =

measure of fit + λ 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)

43

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Lasso regression: L_1 regularized regression

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

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

λ tuning parameter = balance of fit and sparsity

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

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

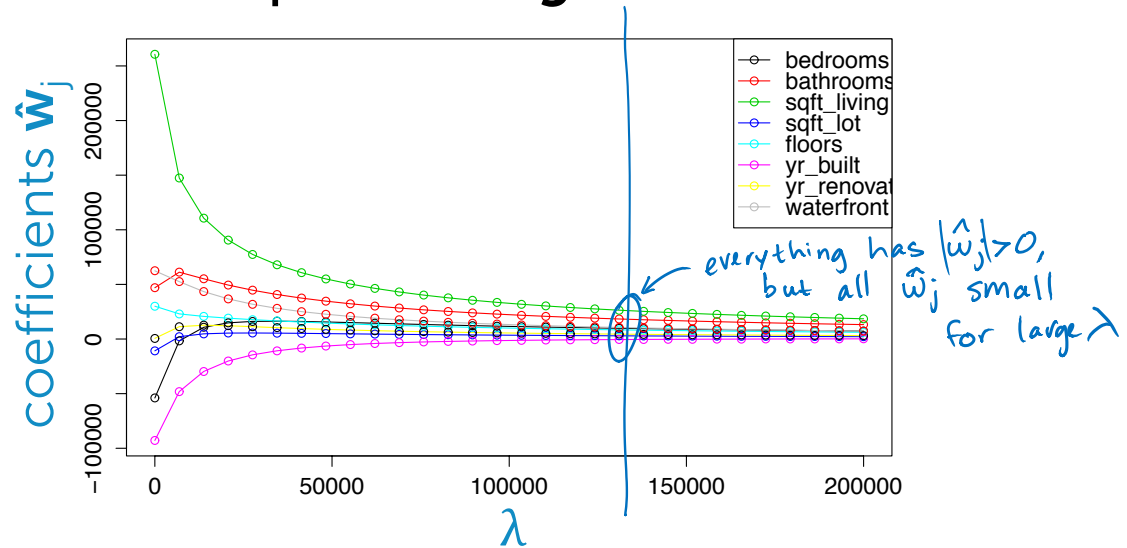
If λ in between: $0 \leq \|\hat{\mathbf{w}}^{\text{lasso}}\|_1 \leq \|\hat{\mathbf{w}}^{\text{LS}}\|_1$

44

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Coefficient path – ridge

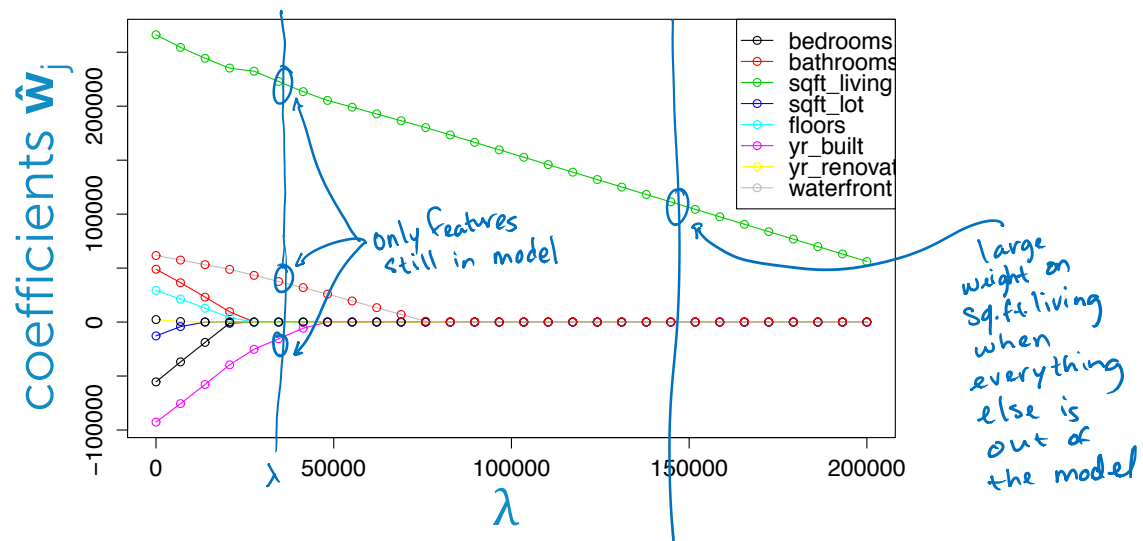


45

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Coefficient path – lasso



46

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

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 ...

47

©2018 Emily Fox

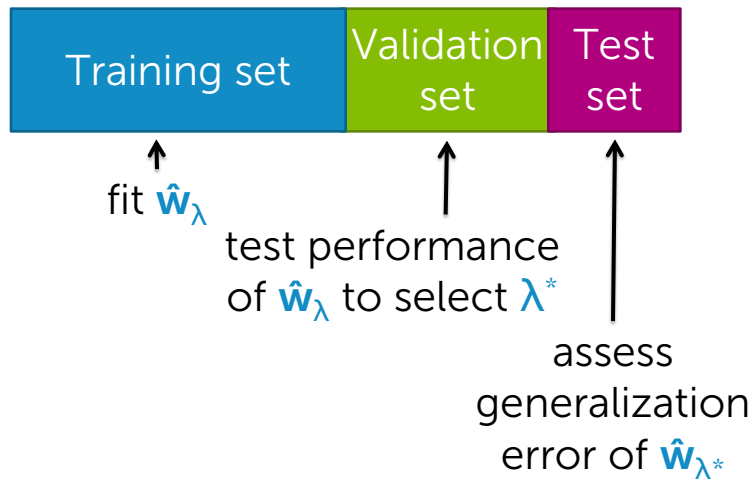
CSE/STAT 416: Intro to Machine Learning

How to choose λ :
Cross validation

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

If sufficient amount of data...



49

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Start with smallish dataset

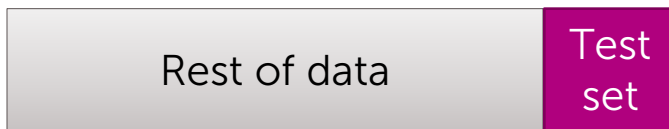
All data

50

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Still form test set and hold out



51

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

How do we use the other data?



use for both training and
validation, but not so naively

52

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Recall naïve approach



↑
small validation set

Is validation set enough to compare performance of $\hat{\mathbf{w}}_\lambda$ across λ values?

No

53

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Choosing the validation set



↑
small validation set

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

Can use any data subset

54

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Choosing the validation set



Which subset should I use?

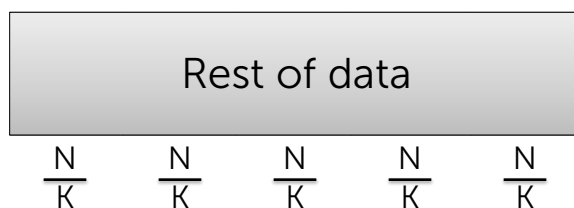
ALL! average performance
over all choices

55

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

K-fold cross validation



Preprocessing: Randomly assign data to K groups

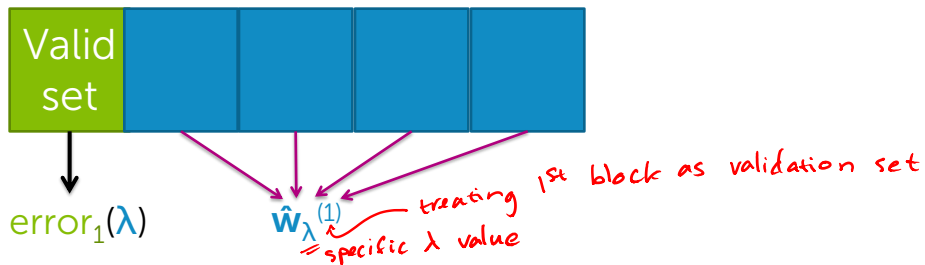
(use same split of data for all other steps)

56

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

K-fold cross validation



For $k=1, \dots, K$

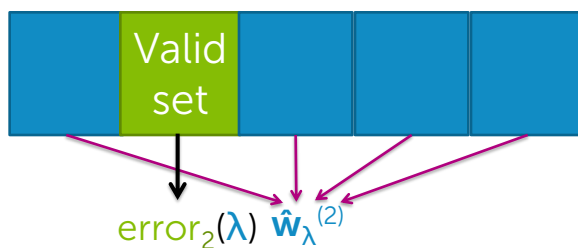
1. Estimate $\hat{w}_\lambda^{(k)}$ on the training blocks
2. Compute error on validation block: $\text{error}_k(\lambda)$

57

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

K-fold cross validation



For $k=1, \dots, K$

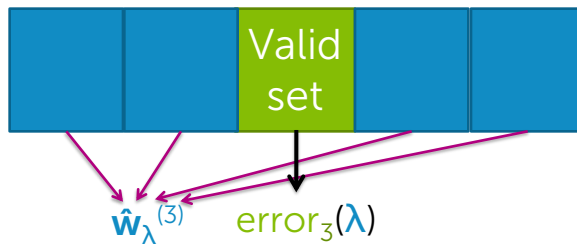
1. Estimate $\hat{w}_\lambda^{(k)}$ on the training blocks
2. Compute error on validation block: $\text{error}_k(\lambda)$

58

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

K-fold cross validation



For $k=1, \dots, K$

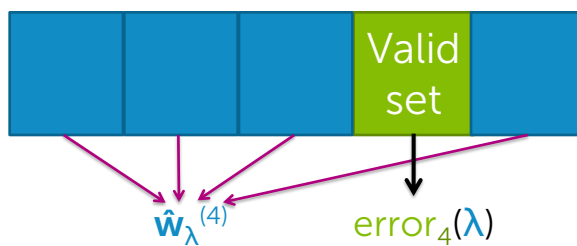
1. Estimate $\hat{w}_\lambda^{(k)}$ on the training blocks
2. Compute error on validation block: $\text{error}_k(\lambda)$

59

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

K-fold cross validation



For $k=1, \dots, K$

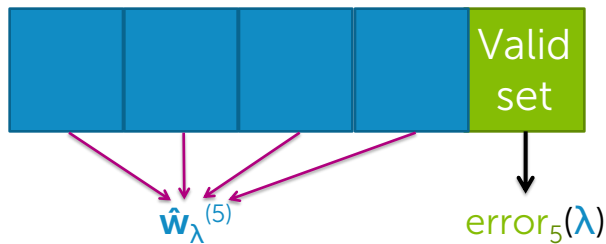
1. Estimate $\hat{w}_\lambda^{(k)}$ on the training blocks
2. Compute error on validation block: $\text{error}_k(\lambda)$

60

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

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

61

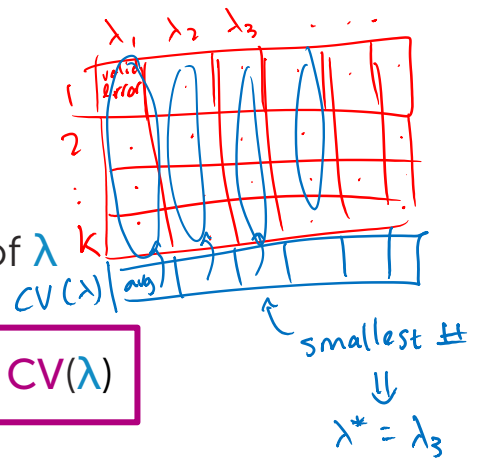
©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

K-fold cross validation



Repeat procedure for each choice of λ



Choose λ^* to minimize $\text{CV}(\lambda)$

62

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

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

63

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

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

c.f., “Machine Learning: A Probabilistic Perspective”, Murphy, 2012 for further discussion

64

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Practical concerns with lasso

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

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

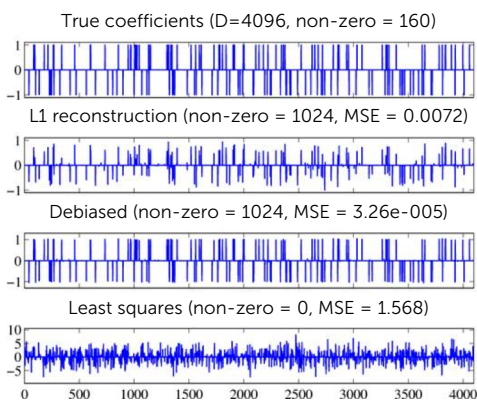


Figure used with permission of Mario Figueiredo
(captions modified to fit course)

66

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

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

Summary for feature selection and lasso regression

Impact of feature selection and lasso

Lasso has changed machine learning, statistics, & electrical engineering

But, for feature selection in general, be **careful about interpreting selected features**

- selection only considers features included
- sensitive to correlations between features
- result depends on algorithm used
- there are theoretical guarantees for lasso under certain conditions

69

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

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
- Implement K-fold cross validation to select lasso tuning parameter λ

70

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning