

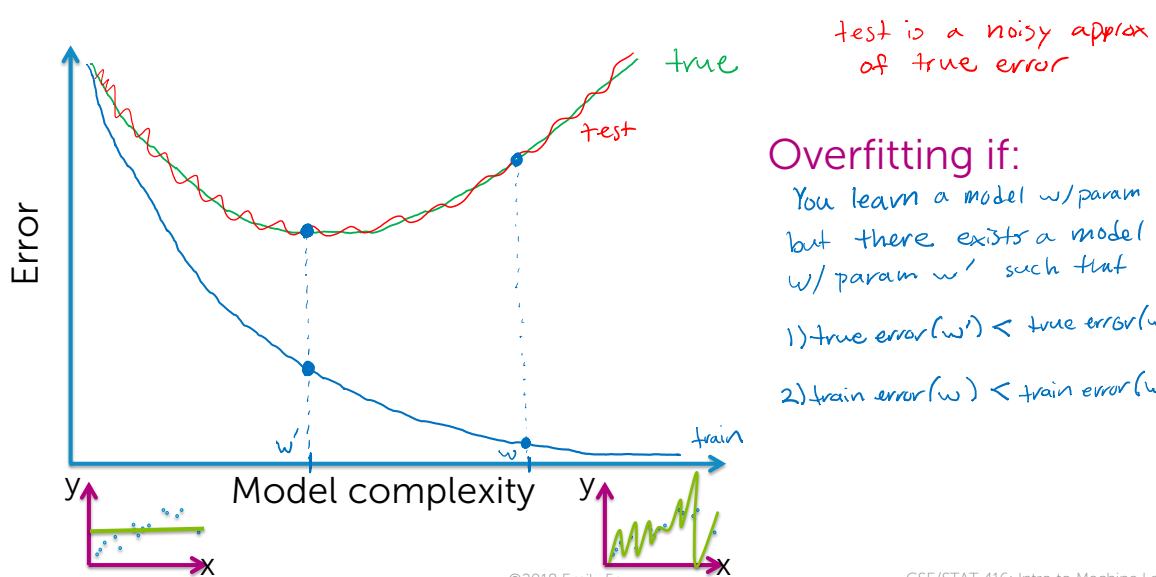
# Ridge Regression:

Regulating overfitting when  
using many features

CSE/STAT 416: Intro to Machine Learning  
Emily Fox  
University of Washington  
April 10, 2018

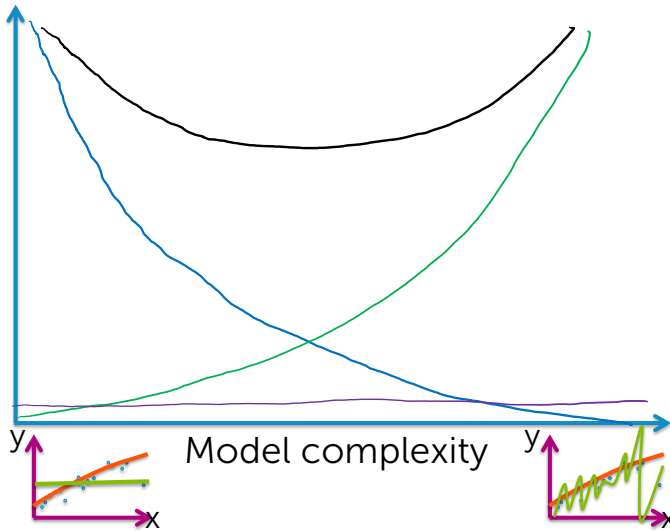
©2018 Emily Fox

## Training, true, & test error vs. model complexity



# Bias-variance tradeoff

$$\text{error} = \text{bias}^2 + \text{variance} + \text{noise}$$



Simple Models

Bias: High

Variance: Low

Complex Models

Bias: Low

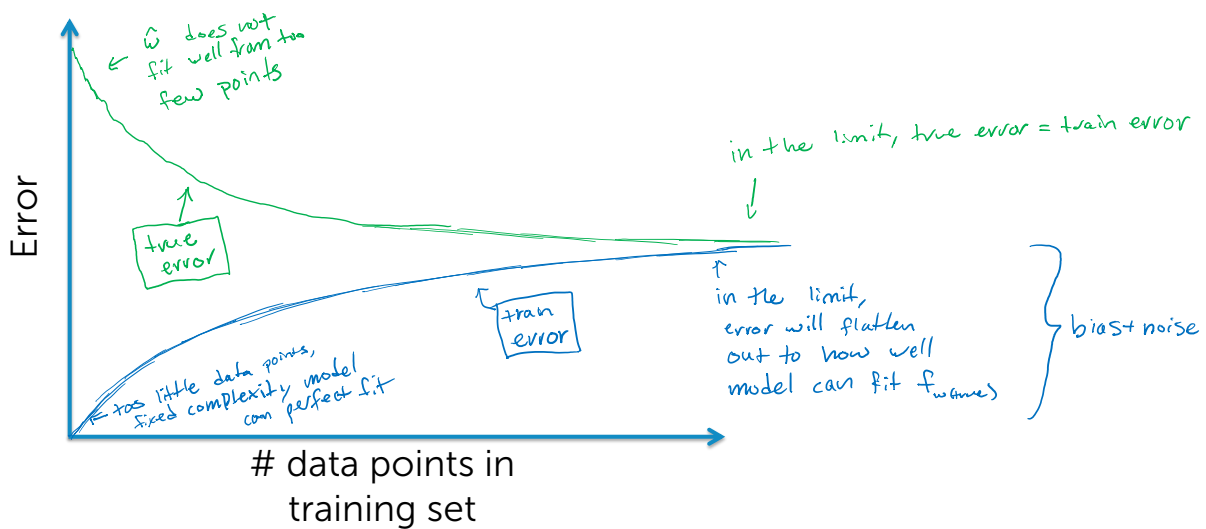
Variance High

3

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

# Error vs. amount of data for fixed model complexity



4

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## Summary of assessing performance

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## What you can do now...

- Describe what a loss function is and give examples
- Contrast training and test error
- Compute training and test error given a loss function
- Discuss issue of assessing performance on training set
- Describe tradeoffs in forming training/test splits
- List and interpret the 3 sources of avg. prediction error
  - Irreducible error, bias, and variance

6

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

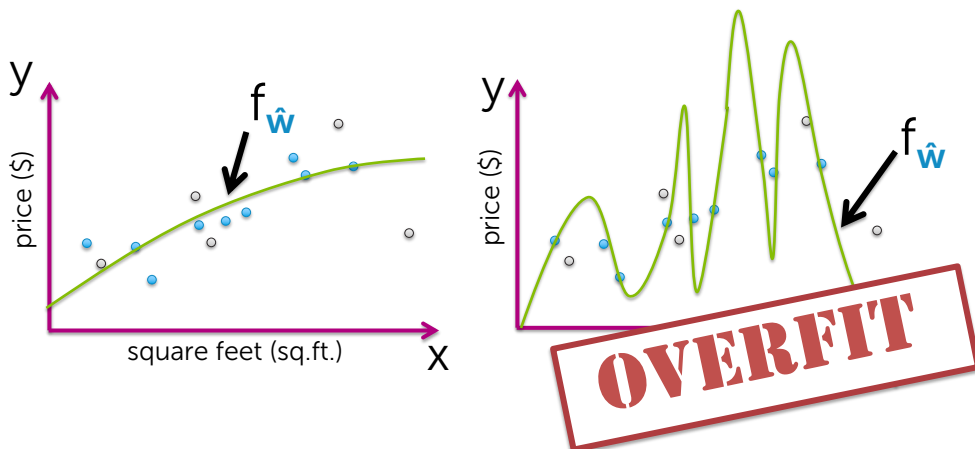
## Overfitting of polynomial regression

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## Flexibility of high-order polynomials

$$y_i = w_0 + w_1 x_i + w_2 x_i^2 + \dots + w_p x_i^p + \varepsilon_i$$



8

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## Symptom of overfitting

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

9

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

Overfitting of linear regression models more generically

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## Overfitting with many features

Not unique to polynomial regression,  
but also if **lots of inputs** ( $d$  large)

Or, generically,  
**lots of features** ( $D$  large)

$$y_i = \sum_{j=0}^D w_j h_j(\mathbf{x}_i) + \varepsilon_i$$

- Square feet
- # bathrooms
- # bedrooms
- Lot size
- Year built
- ...

11

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

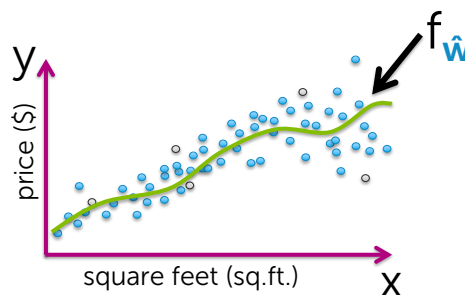
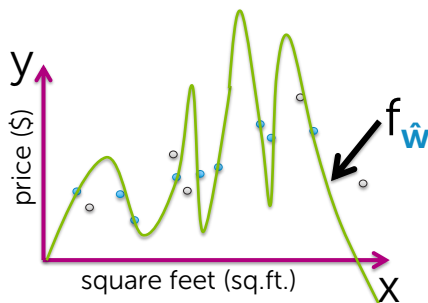
## How does # of observations influence overfitting?

Few observations ( $N$  small)

→ rapidly overfit as model complexity increases

Many observations ( $N$  very large)

→ harder to overfit



12

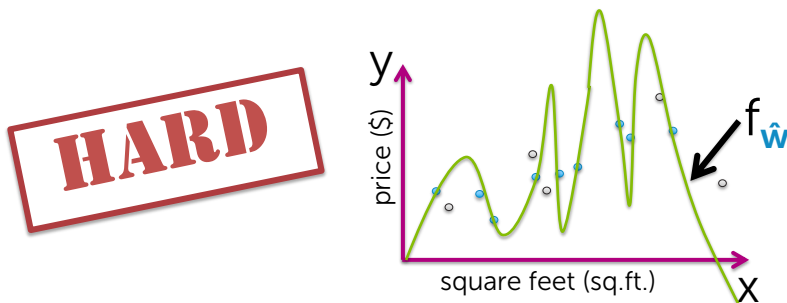
©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## How does # of inputs influence overfitting?

**1 input** (e.g., sq.ft.):

Data must include representative examples of all possible (sq.ft., \$) pairs to avoid overfitting



13

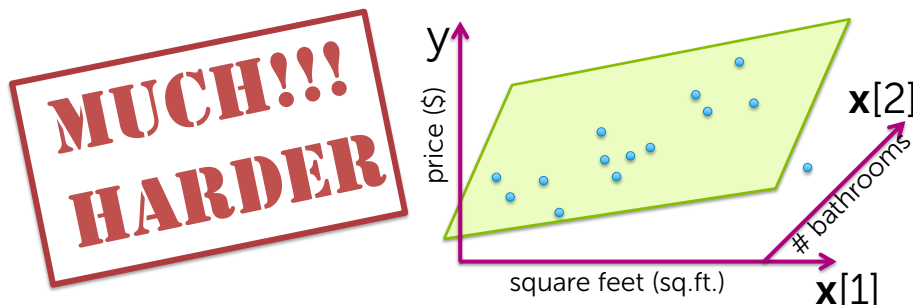
©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## How does # of inputs influence overfitting?

**d inputs** (e.g., sq.ft., #bath, #bed, lot size, year,...):

Data must include examples of all possible (sq.ft., #bath, #bed, lot size, year,..., \$) combos to avoid overfitting



14

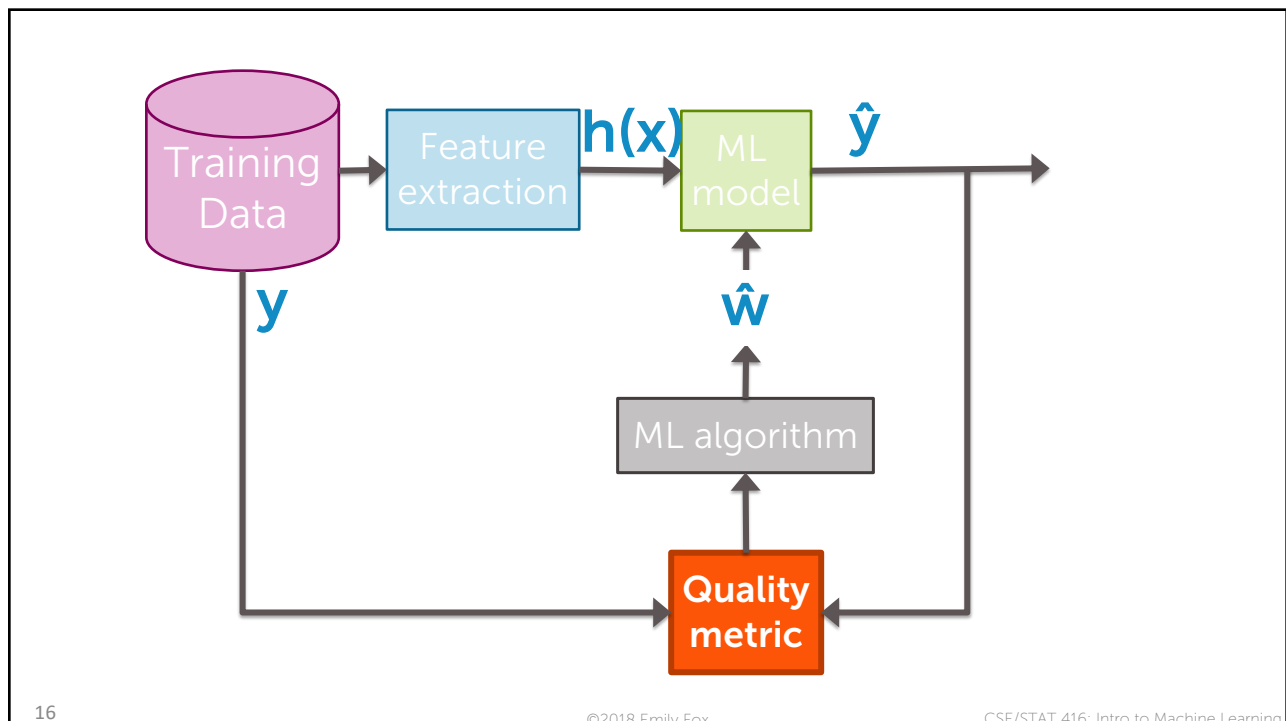
©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## Adding term to cost-of-fit to prefer small coefficients

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning



16

©2018 Emily Fox

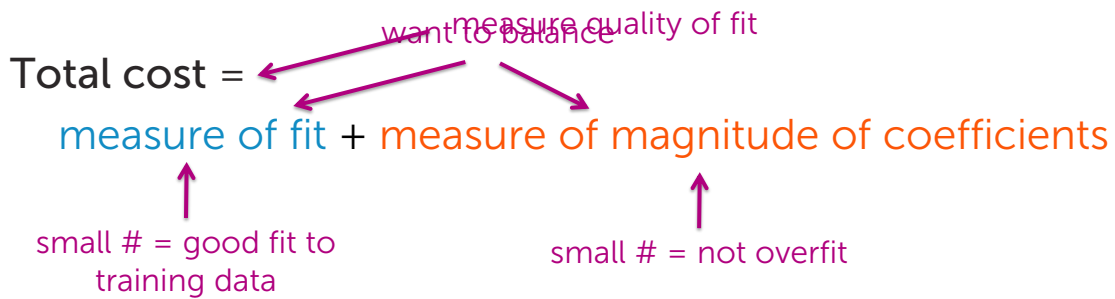
CSE/STAT 416: Intro to Machine Learning



## Desired total cost format

Want to balance:

- i. How well function fits data
- ii. Magnitude of coefficients

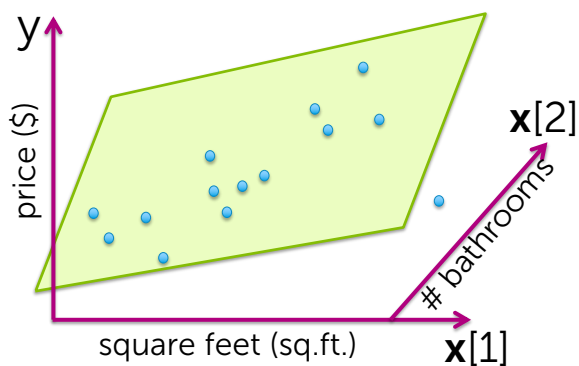


17

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## Measure of fit to training data



$$\text{RSS}(\mathbf{w}) = \sum_{i=1}^N (y_i - \mathbf{h}(\mathbf{x}_i)^T \mathbf{w})^2$$

$$= \sum_{i=1}^N (y_i - \hat{y}_i(\mathbf{w}))^2$$

small RSS  $\rightarrow$  model fitting training data well

18

©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?  $\sum_{j=0}^D |w_j| \triangleq \|w\|_1$   $w = [w_0 w_1 \dots w_D]$
- Sum of squares ( $L_2$  norm)  $\sum_{j=0}^D w_j^2 \triangleq \|w\|_2^2$  ←  $L_1$  norm... discuss next lecture  
← focus of this lecture

19

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## Consider specific total cost

Total cost =

measure of fit + measure of magnitude of coefficients

20

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## Consider specific total cost

Total cost =

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

21

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## Consider resulting objective

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

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

$\lambda$  tuning parameter = balance of fit and magnitude

If  $\lambda=0$ :  
reduces to  $\min \text{RSS}(\mathbf{w})$ , as before (old soln)  
 $\rightarrow \hat{\mathbf{w}}^{\text{LS}}$  (least squares)

If  $\lambda=\infty$ :  
For solns  $\hat{\mathbf{w}} \neq 0$ , then total cost  $\rightarrow \infty$   
If  $\hat{\mathbf{w}} = 0$ , then total cost =  $\text{RSS}(0)$   $\rightarrow \hat{\mathbf{w}} = 0$

If  $\lambda$  in between: Then  $0 \leq \|\hat{\mathbf{w}}\|_2^2 \leq \|\hat{\mathbf{w}}^{\text{LS}}\|_2^2$

22

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## 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)

23

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## Bias-variance tradeoff

Large  $\lambda$ :

high bias, low variance

(e.g.,  $\hat{\mathbf{w}} = 0$  for  $\lambda = \infty$ )

In essence,  $\lambda$   
controls model  
complexity

Small  $\lambda$ :

low bias, high variance

(e.g., standard least squares (RSS) fit of  
high-order polynomial for  $\lambda = 0$ )

24

©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 ridge regression?

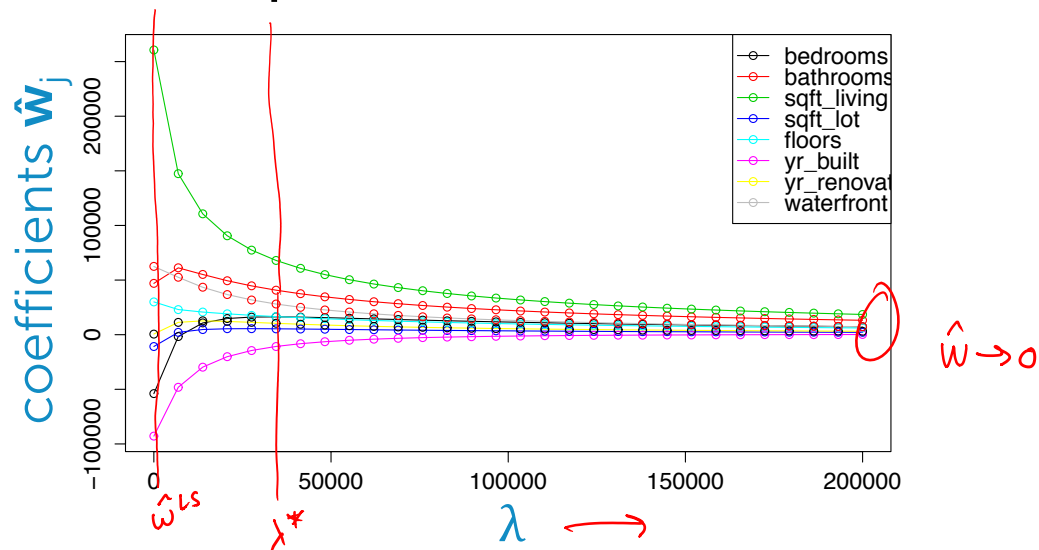
Will consider a few settings of  $\lambda$  ...

25

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## Coefficient path



26

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## How to choose $\lambda$

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## The regression/ML workflow

1. Model selection  
Need to **choose tuning parameters  $\lambda$**  controlling model complexity
2. Model assessment  
Having selected a model, **assess generalization error**

28

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## Hypothetical implementation

Training set

Test set

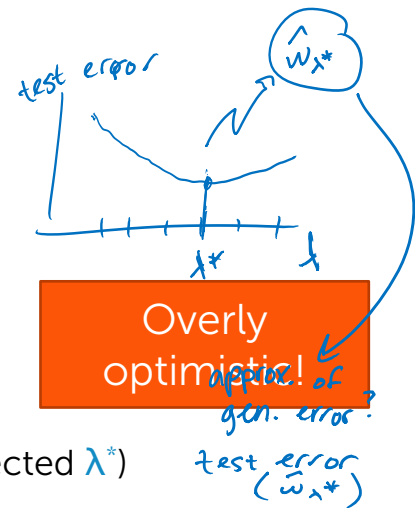
### 1. Model selection

For each considered  $\lambda$  :

- i. Estimate parameters  $\hat{\mathbf{w}}_\lambda$  on training data
- ii. Assess performance of  $\hat{\mathbf{w}}_\lambda$  on test data
- iii. Choose  $\lambda^*$  to be  $\lambda$  with lowest test error

### 2. Model assessment

Compute test error of  $\hat{\mathbf{w}}_{\lambda^*}$  (fitted model for selected  $\lambda^*$ ) to approx. true error



29

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## Hypothetical implementation

Training set

Test set

**Issue:** Just like fitting  $\hat{\mathbf{w}}$  and assessing its performance both on training data

- $\lambda^*$  was selected to minimize test error (i.e.,  $\lambda^*$  was fit on test data)
- If test data is not representative of the whole world, then  $\hat{\mathbf{w}}_{\lambda^*}$  will typically perform worse than test error indicates

30

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## Practical implementation



Solution: Create two “test” sets!

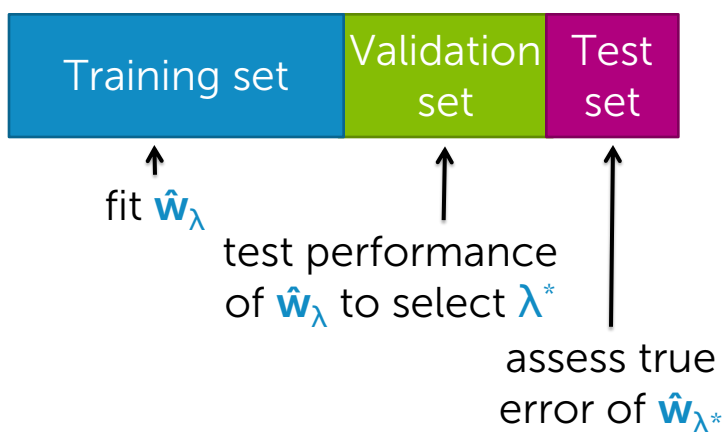
1. Select  $\lambda^*$  such that  $\hat{w}_{\lambda^*}$  minimizes error on validation set
2. Approximate true error of  $\hat{w}_{\lambda^*}$  using test set

31

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## Practical implementation



32

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning



## Typical splits

Training set	Validation set	Test set
80%	10%	10%
50%	25%	25%

How to handle the intercept

**PRACTICALITIES**

## Recall multiple regression model

Model:

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

$$= \sum_{j=0}^D w_j h_j(\mathbf{x}_i) + \varepsilon_i$$

feature 1 =  $h_0(\mathbf{x})$ ... often 1 (constant)

feature 2 =  $h_1(\mathbf{x})$ ... e.g.,  $\mathbf{x}[1]$

feature 3 =  $h_2(\mathbf{x})$ ... e.g.,  $\mathbf{x}[2]$

...

feature  $D+1$  =  $h_D(\mathbf{x})$ ... e.g.,  $\mathbf{x}[d]$

Assume our 1st feature is 1 (include intercept)

35

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## Do we penalize intercept?

Standard ridge regression cost:

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

strength of penalty

Encourages intercept  $w_0$  to also be small

intercept  
[ $w_0, w_1, \dots, w_D$ ]

Do we want a small intercept?

Conceptually, not indicative of overfitting...

36

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## Option 1: Don't penalize intercept

Modified ridge regression cost:

$$\text{RSS}(w_0, \mathbf{w}_{\text{rest}}) + \lambda \|\mathbf{w}_{\text{rest}}\|_2^2$$


  
 $[w_1 \dots w_D]$

37

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## Option 2: Center data first

If data are first **centered about 0**, then favoring small intercept not so worrisome

**Step 1:** Transform  $y$  to have 0 mean

**Step 2:** Run ridge regression as normal  
(closed-form or gradient algorithms)

38

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## Feature normalization

**PRACTICALITIES**

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## 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:  $Z_j$

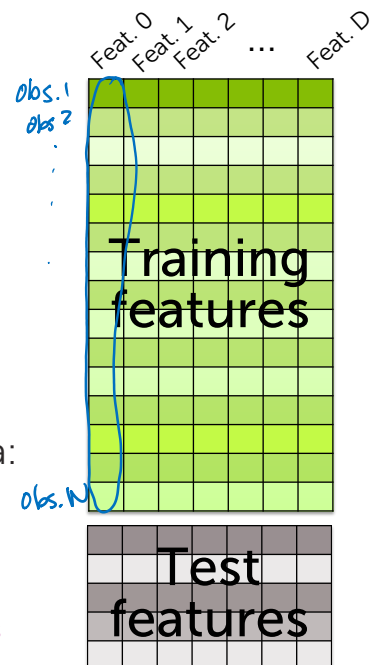
Apply same training scale factors to test data:

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

apply to test point

Normalizer:  $Z_j$

summing over training points



40

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

## Summary for ridge regression

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning

### What you can do now...

- Describe what happens to magnitude of estimated coefficients when model is overfit
- Motivate form of ridge regression cost function
- Describe what happens to estimated coefficients of ridge regression as tuning parameter  $\lambda$  is varied
- Interpret coefficient path plot
- Use a validation set to select the ridge regression tuning parameter  $\lambda$
- Handle intercept and scale of features with care

42

©2018 Emily Fox

CSE/STAT 416: Intro to Machine Learning