# Ridge Regression:
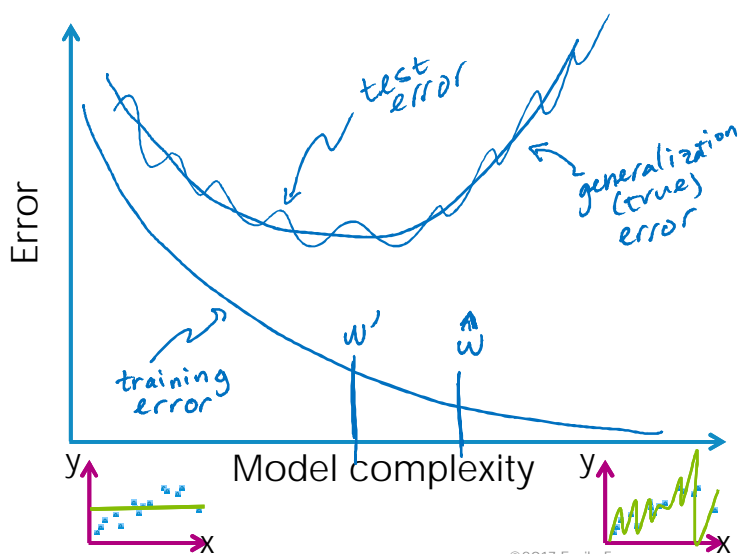## Regulating overfitting when using many features

CSE 446: Machine Learning
Emily Fox
University of Washington
January 13, 2017

---

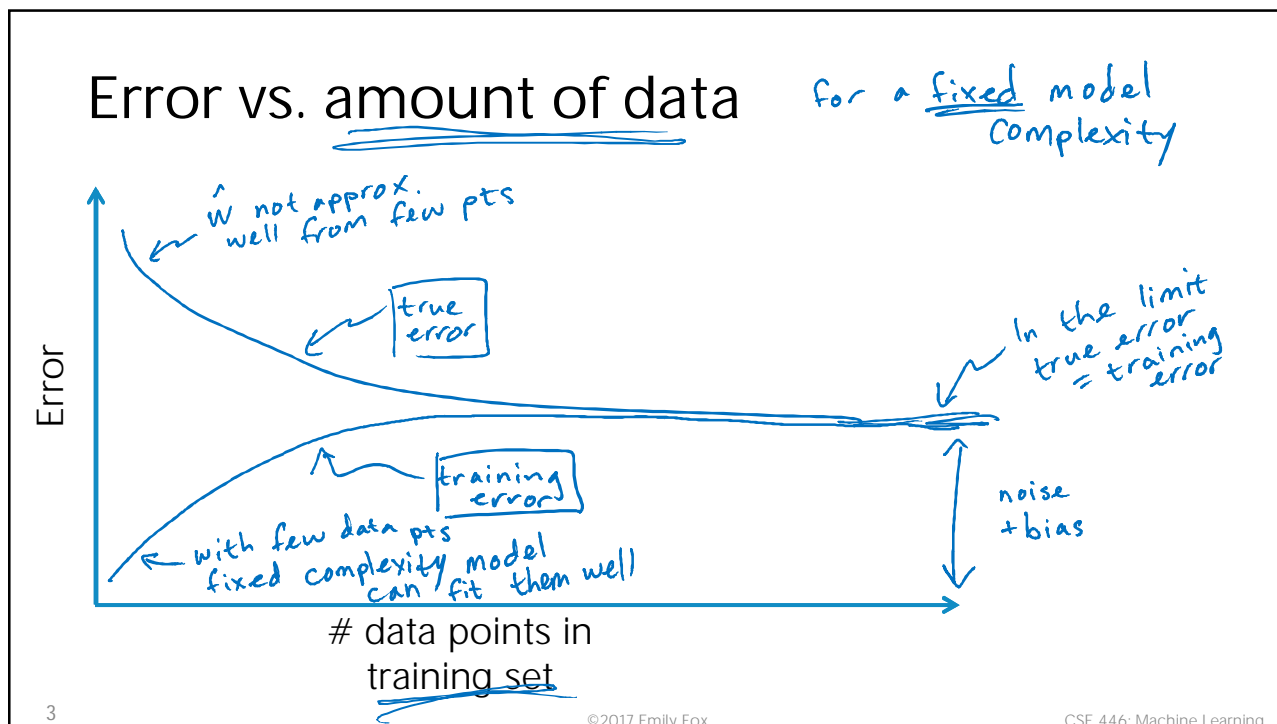# Training, true, & test error vs. model complexity



Overfitting if:

If there exists a model with estimated params $w'$ such that

① training error $(\hat{w})$ < training error $(w')$

② true error $(\hat{w})$ > true error $(w')$

2

CSE 446: Machine Learning

# Error vs. amount of data

*for a fixed model Complexity*



$\hat{w}$ not approx. well from few pts

true error

Error

In the limit true error = training error

training error

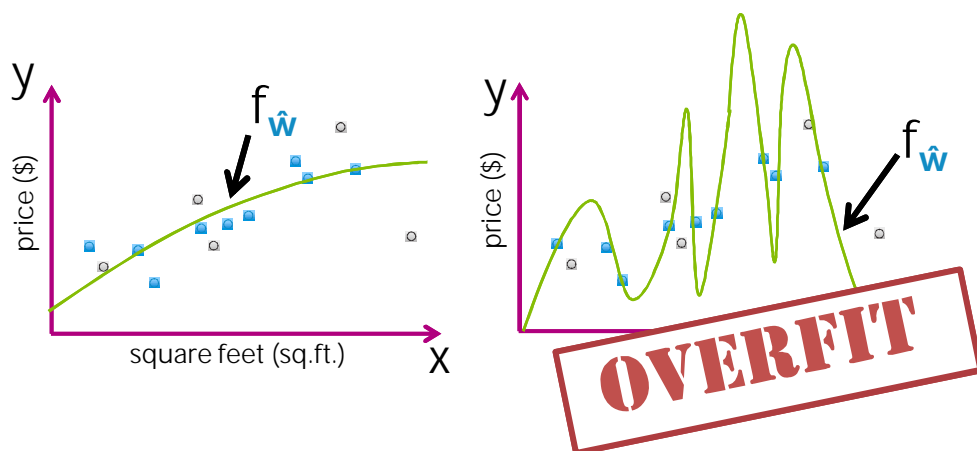with few data pts fixed complexity model can fit them well

noise + bias

# data points in training set

# Overfitting of polynomial regression

# Flexibility of high-order polynomials

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



OVERFIT

©2017 Emily Fox                                CSE 446: Machine Learning

# Symptom of overfitting

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

©2017 Emily Fox                                CSE 446: Machine Learning

# Overfitting of linear regression models more generically

# 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
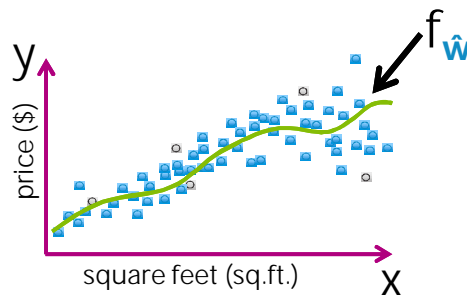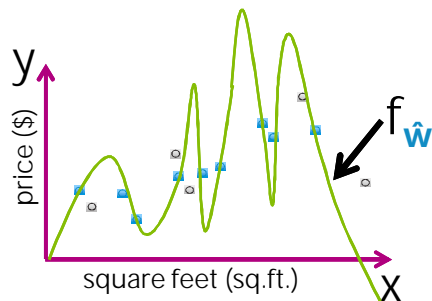- Lot size
- Year built
- ...

CSE 446: Machine Learning

# How does # of observations influence overfitting?

Few **observations** (N small)
→ rapidly overfit as model complexity increases
Many **observations** (N very large)
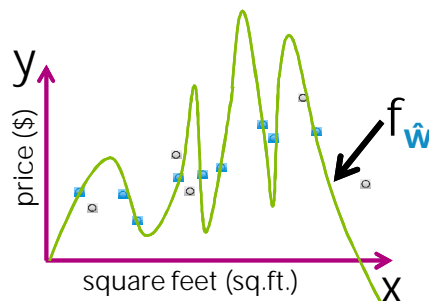→ harder to overfit

©2017 Emily Fox
CSE 446: Machine Learning

# How does # of inputs influence overfitting?

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

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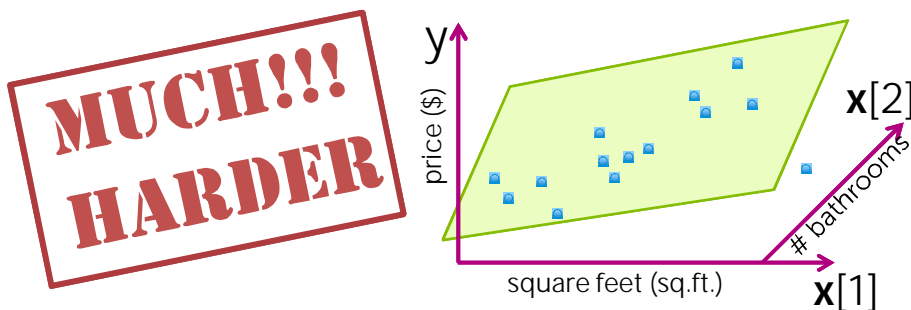all possible (sq.ft., $) pairs to avoid overfitting

HARD

©2017 Emily Fox
CSE 446: 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

MUCH!!!
HARDER



11              ©2017 Emily Fox                    CSE 446: Machine Learning

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

# Desired total cost format

Want to balance:

i. How well function fits data

ii. Magnitude of coefficients

want to balance measure quality of fit
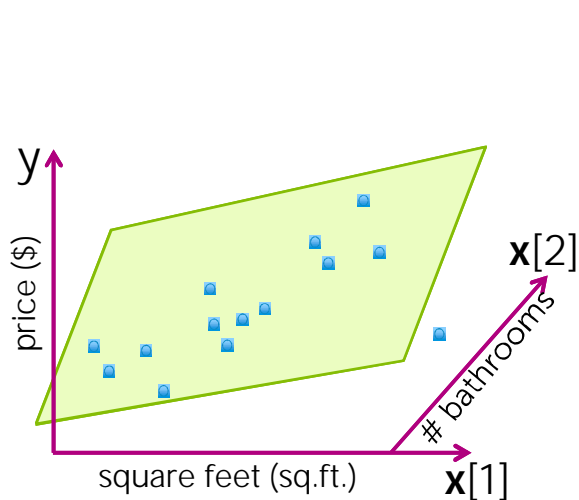
Total cost =

measure of fit + measure of magnitude of coefficients

small # = good fit to training data

small # = not overfit

©2017 Emily Fox · CSE 446: Machine Learning

---

# Measure of fit to training data



y

price ($)

**x**[2]

# bathrooms

square feet (sq.ft.)

**x**[1]

pred. value using w

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

small RSS → model fitting training data well

©2017 Emily Fox · CSE 446: 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 = $ small # ✕

- Sum of absolute value?

  $$\sum_{j=0}^{D} |w_j| \triangleq \|w\|_1$$     $L_1$ norm... discuss next lecture

- Sum of squares ($L_2$ norm)

  $$\sum_{j=0}^{D} w_j^2 \triangleq \|w\|_2^2$$     $L_2$ norm... focus of this lecture

15

©2017 Emily Fox

CSE 446: Machine Learning

# Consider specific total cost

Total cost =
    measure of fit + measure of magnitude of coefficients

16

©2017 Emily Fox

CSE 446: Machine Learning

# Consider specific total cost

Total cost =
    <span style="color:blue">measure of fit</span> + <span style="color:orange">measure of magnitude of coefficients</span>

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

©2017 Emily Fox                                    CSE 446: 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

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

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

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

©2017 Emily Fox                                    CSE 446: 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)

©2017 Emily Fox CSE 446: 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 $\lambda = 0$)

©2017 Emily Fox CSE 446: Machine Learning

# Revisit polynomial fit demo

What happens if we refit our high-order polynomial, but now using **ridge regression**?
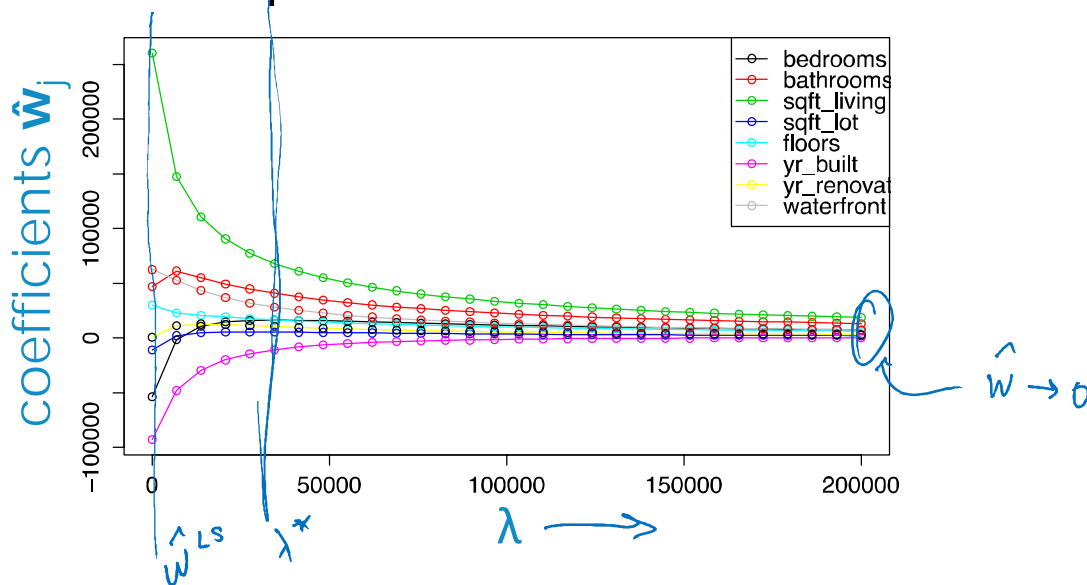
Will consider a few settings of $\lambda$ ...

21          ©2017 Emily Fox          CSE 446: Machine Learning

# Coefficient path



22          ©2017 Emily Fox          CSE 446: Machine Learning

Fitting the ridge regression model
(for given λ value)

**Step 1:**
Rewrite total cost in matrix notation

# Recall matrix form of RSS

Model for all N observations together



$$\mathbf{y} = \mathbf{H}\,\mathbf{w} + \boldsymbol{\varepsilon}$$

©2017 Emily Fox                                    CSE 446: Machine Learning

# Recall matrix form of RSS

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

$$= (\mathbf{y}-\mathbf{H}\mathbf{w})^\mathsf{T}(\mathbf{y}-\mathbf{H}\mathbf{w})$$

©2017 Emily Fox                                    CSE 446: Machine Learning

## Rewrite magnitude of coefficients in vector notation

$$||\mathbf{w}||_2^2 = w_0{}^2 + w_1{}^2 + w_2{}^2 + ... + w_D{}^2$$

$$= \begin{bmatrix} & & & & & \\ w_0 & w_1 & \cdots & & w_D \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ \\ \\ \\ w_D \end{bmatrix}$$

$$= \mathbf{w}^\top \mathbf{w}$$

27

©2017 Emily Fox

CSE 446: Machine Learning

## Putting it all together

In matrix form, ridge regression cost is:

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

$$= (\mathbf{y} - \mathbf{H}\mathbf{w})^\top (\mathbf{y} - \mathbf{H}\mathbf{w}) + \lambda \mathbf{w}^\top \mathbf{w}$$

28

©2017 Emily Fox

CSE 446: Machine Learning

14

**Step 2:**
Compute the gradient

---

# Gradient of ridge regression cost

$$\|w\|_2$$
$$= \sqrt{w_1^2 + \ldots + w_p^2}$$
$$= \sqrt{w^\top w}$$

$$\nabla\left[\text{RSS}(w) + \lambda\|w\|_2^2\right] = \nabla\left[(y-Hw)^\top(y-Hw) + \lambda w^\top w\right]$$

$$= \underbrace{\nabla\left[(y-Hw)^\top(y-Hw)\right]}_{-2H^\top(y-Hw)} + \lambda\underbrace{\nabla[w^\top w]}_{2w}$$

**Why?** By analogy to 1d case...

$w^\top w$ analogous to $w^2$ and derivative of $w^2 = 2w$

CSE 446: Machine Learning

## Step 3, Approach 1:
Set the gradient = 0

# Ridge closed-form solution

3D plot of RSS with tangent plane at minimum

$$\nabla \text{cost}(\mathbf{w}) = -2\mathbf{H}^T(\mathbf{y}-\mathbf{Hw}) + 2\lambda\mathbf{I}\mathbf{w} = 0$$
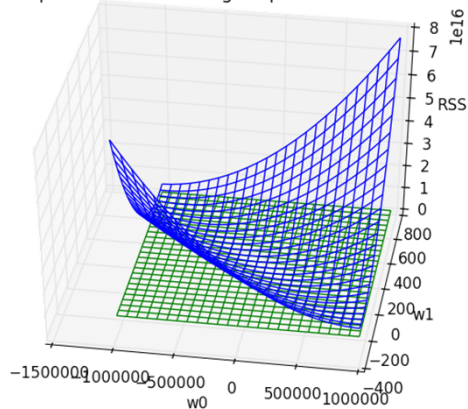
add

Solve for **w**:

$$-\mathbf{H}^T y + \mathbf{H}^T \mathbf{H} \hat{w} + \lambda \mathbf{I} \hat{w} = 0$$

$$(\mathbf{H}^T \mathbf{H} + \lambda \mathbf{I}) \hat{w} = \mathbf{H}^T y$$

$$\hat{w}^{ridge} = (\mathbf{H}^T \mathbf{H} + \lambda \mathbf{I})^{-1} \mathbf{H}^T y$$

32

CSE 446: Machine Learning

# Interpreting ridge closed-form solution

3D plot of RSS with tangent plane at minimum

$$\hat{\mathbf{w}} = (\mathbf{H}^T\mathbf{H} + \lambda\mathbf{I})^{-1}\mathbf{H}^T\mathbf{y}$$

If $\lambda = 0$:    $\hat{w}^{ridge} = \left(H^T H\right)^{-1} H^T y = \hat{w}^{LS}$

old soln ✓

If $\lambda = \infty$:    $\hat{w}^{ridge} = 0$ ← because it's like dividing by ∞ ✓

33    ©2017 Emily Fox    CSE 446: Machine Learning

---

# Recall discussion on previous closed-form solution

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

# feat D

# obs N

Invertible if:
   In general,
   (# linearly independent obs)
   N > D

Complexity of inverse:
   $O(D^3)$
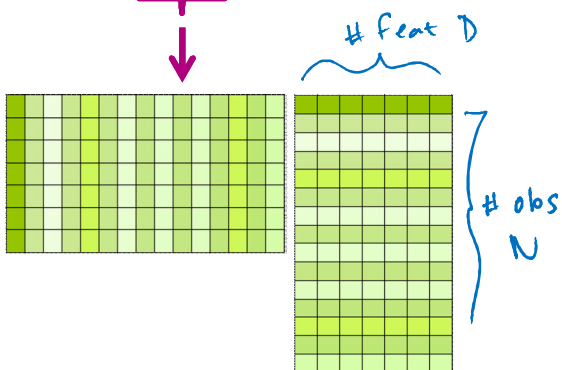
34    ©2017 Emily Fox    CSE 446: Machine Learning

17

# Discussion of ridge closed-form solution

$$\hat{\mathbf{w}} = (\mathbf{H}^T\mathbf{H} + \lambda\mathbf{I})^{-1}\mathbf{H}^T\mathbf{y}$$

*λI is making $H^TH + \lambda I$ more "regular"*
*→ "regularization"*

Invertible if:
Always if λ>0,
even if N < D

Complexity of inverse:
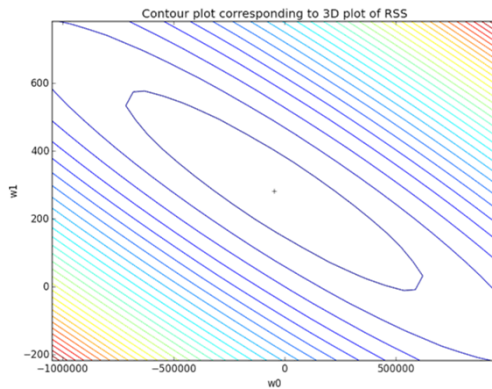$O(D^3)$...
big for large D!

35

©2017 Emily Fox

CSE 446: Machine Learning

---

## Step 3, Approach 2:
Gradient descent

# Elementwise ridge regression gradient descent algorithm

$$\nabla \text{cost}(\mathbf{w}) = -2\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{w}) + 2\lambda\mathbf{w}$$

Contour plot corresponding to 3D plot of RSS

Update to $j^{th}$ feature weight:

$$w_j^{(t+1)} \leftarrow w_j^{(t)} - \eta \,*$$

as before $\left[ -2\sum_{i=1}^{N} h_j(\mathbf{x}_i)(y_i - \hat{y}_i(\mathbf{w}^{(t)})) \right.$

new term $\left. +2\lambda w_j^{(t)} \right]$

37
©2017 Emily Fox
CSE 446: Machine Learning

# Recall previous algorithm

init $\mathbf{w}^{(1)} = 0$ (or randomly, or smartly), $t=1$

**while** $|| \nabla \text{RSS}(\mathbf{w}^{(t)})|| > \varepsilon$

    **for** j=0,...,D

        partial[j] $= -2\sum_{i=1}^{N} h_j(\mathbf{x}_i)(y_i - \hat{y}_i(\mathbf{w}^{(t)}))$

        $w_j^{(t+1)} \leftarrow w_j^{(t)} - \eta$ partial[j]

    $t \leftarrow t + 1$

38
©2017 Emily Fox
CSE 446: Machine Learning

## Summary of ridge regression algorithm

init $\mathbf{w}^{(1)}$=0 (or randomly, or smartly), $t$=1

**while** $||\nabla\text{RSS}(\mathbf{w}^{(t)})|| > \varepsilon$

    **for** $j$=0,...,D

        $\text{partial}[j] = -2\sum_{i=1}^{N}h_j(\mathbf{x}_i)(y_i - \hat{y}_i(\mathbf{w}^{(t)}))$

        $w_j^{(t+1)} \leftarrow (1\text{-}2\eta\lambda)w_j^{(t)} - \eta\ \text{partial}[j]$

    $t \leftarrow t + 1$

39      ©2017 Emily Fox      CSE 446: Machine Learning

## How to choose λ

©2017 Emily Fox

# 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

©2017 Emily Fox CSE 446: 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

Overly optimistic!

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

©2017 Emily Fox CSE 446: 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

©2017 Emily Fox  CSE 446: Machine Learning

# Practical implementation

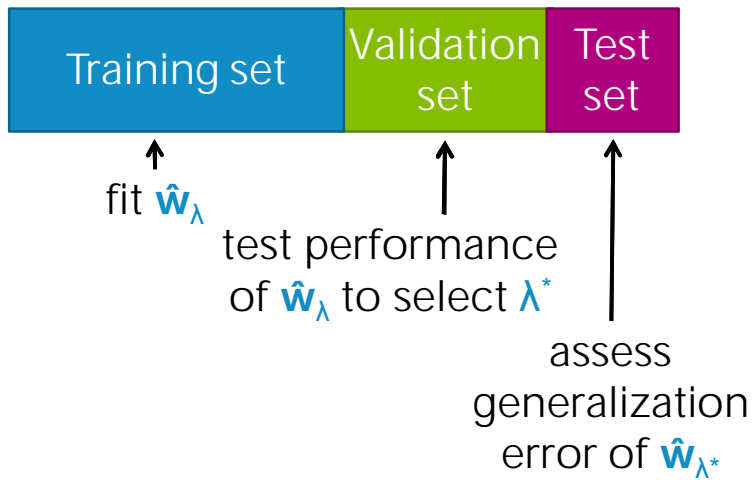| Training set | Validation set | Test set |
|---|---|---|

Solution: Create two "test" sets!

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

©2017 Emily Fox  CSE 446: Machine Learning

この上部の日付

# Practical implementation

| Training set | Validation set | Test set |
|---|---|---|

fit $\hat{\mathbf{w}}_\lambda$

test performance of $\hat{\mathbf{w}}_\lambda$ to select $\lambda^*$

assess generalization error of $\hat{\mathbf{w}}_{\lambda^*}$

©2017 Emily Fox
CSE 446: Machine Learning

# Typical splits

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

©2017 Emily Fox
CSE 446: Machine Learning

# How to handle the intercept

OPTIONAL

©2017 Emily Fox

---

# Recall multiple regression model

Model:

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

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

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., $\mathbf{x}[d]$

48   ©2017 Emily Fox   CSE 446: Machine Learning

## If constant feature…

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

In matrix notation for N observations:

$$\mathbf{y} = \mathbf{H}\,\mathbf{w} + \boldsymbol{\varepsilon}$$

$w_0$

©2017 Emily Fox — CSE 446: 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

Do we want a small intercept?
Conceptually, not indicative of overfitting…

©2017 Emily Fox — CSE 446: Machine Learning

## Option 1: Don't penalize intercept

Modified ridge regression cost:

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

How to implement this in practice?

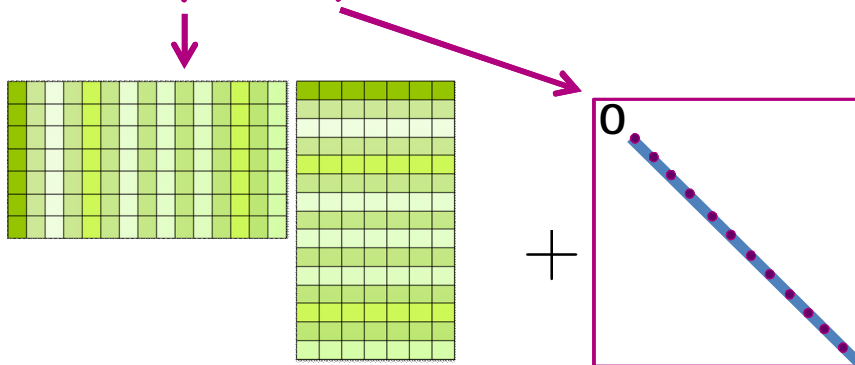©2017 Emily Fox    CSE 446: Machine Learning

## Option 1: Don't penalize intercept
## — Closed-form solution —

$$\hat{\mathbf{w}} = (\mathbf{H}^T\mathbf{H} + \lambda \mathbf{I}^{mod})^{-1} \mathbf{H}^T\mathbf{y}$$



©2017 Emily Fox    CSE 446: Machine Learning

## Option 1: Don't penalize intercept — Gradient descent algorithm —

**while** $|| \nabla \text{RSS}(\mathbf{w}^{(t)})|| > \varepsilon$

    **for** j=0,...,D

        partial[j] $= -2\sum_{i=1}^{N} h_j(\mathbf{x}_i)(y_i - \hat{y}_i(\mathbf{w}^{(t)}))$

        if j==0

            $w_0^{(t+1)} \leftarrow w_0^{(t)} - \eta$ partial[j]

        else

            $w_j^{(t+1)} \leftarrow (1-2\eta\lambda)w_j^{(t)} - \eta$ partial[j]

    t $\leftarrow$ t + 1

53          ©2017 Emily Fox          CSE 446: 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)

54          ©2017 Emily Fox          CSE 446: Machine Learning

# Summary for
# ridge regression

©2017 Emily Fox

# 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
- Estimate ridge regression parameters:
  - In closed form
  - Using an iterative gradient descent algorithm
- Use a validation set to select the ridge regression tuning parameter $\lambda$