

# Lasso Regression:

## Regularization for feature selection

CSE 446: Machine Learning  
Emily Fox  
University of Washington  
January 18, 2017

1

©2017 Emily Fox

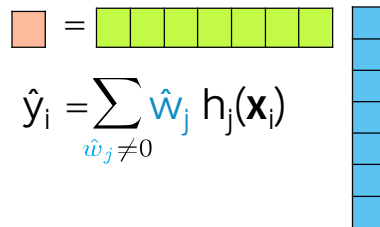
Feature selection task

# 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

← many zeros



## Interpretability:

- Which features are relevant for prediction?

# 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        | ⋮                |

## Option 1: All subsets or greedy variants

### 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 ... 0 0 0]	} $2^{D+1}$	$2^8 = 256$
$y_i = w_0 h_0(\mathbf{x}_i) + \epsilon_i$	[1 0 0 ... 0 0 0]		$2^{30} = 1,073,741,824$
$y_i = w_1 h_1(\mathbf{x}_i) + \epsilon_i$	[0 1 0 ... 0 0 0]		$2^{1000} = 1.071509 \times 10^{301}$
$\vdots$	$\vdots$		$2^{100B} = \text{HUGE!!!!!!}$
$y_i = w_0 h_0(\mathbf{x}_i) + w_1 h_1(\mathbf{x}_i) + \epsilon_i$	[1 1 0 ... 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 ... 1 1 1]		

*feature 0*  
*feature 1*     ...     *feature D*  
 $2 \quad 2 \quad \dots \quad 2$

Typically,  
computationally  
infeasible

## Choosing model complexity?

Option 1: Assess on validation set

Option 2: Cross validation

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

7

©2017 Emily Fox

CSF 446: 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.

8

©2017 Emily Fox

CSF 446: Machine Learning

## Option 2: Regularize

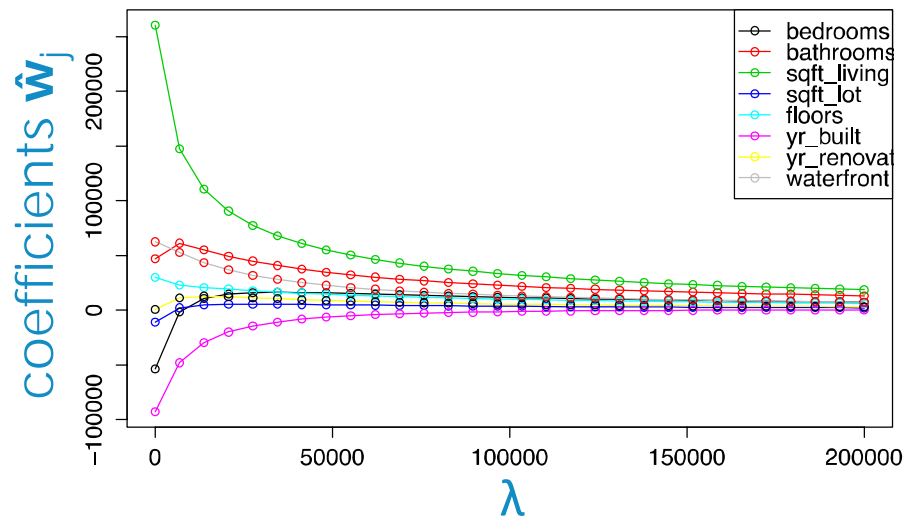
### Ridge regression: $L_2$ regularized regression

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

## Coefficient path – ridge



11

©2017 Emily Fox

CSF 446: 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

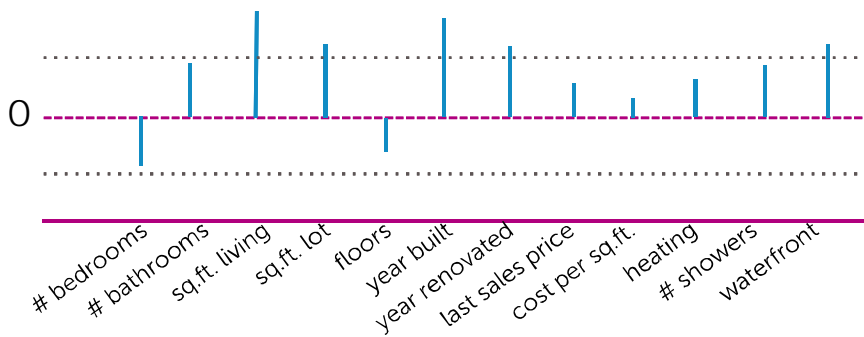
12

©2017 Emily Fox

CSF 446: Machine Learning

# Thresholding ridge coefficients?

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



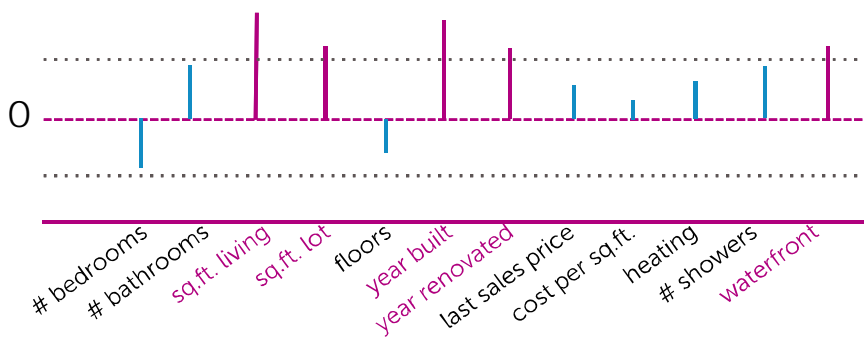
13

©2017 Emily Fox

CSF 446: Machine Learning

# Thresholding ridge coefficients?

Selected features for a given threshold value



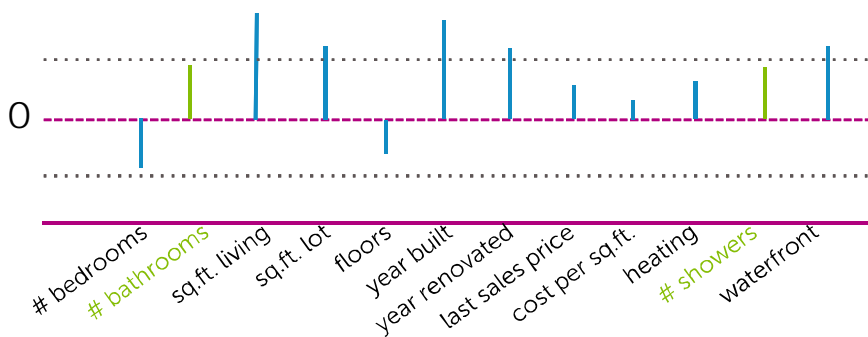
14

©2017 Emily Fox

CSF 446: Machine Learning

# Thresholding ridge coefficients?

Let's look at two related features...



Nothing measuring bathrooms was included!

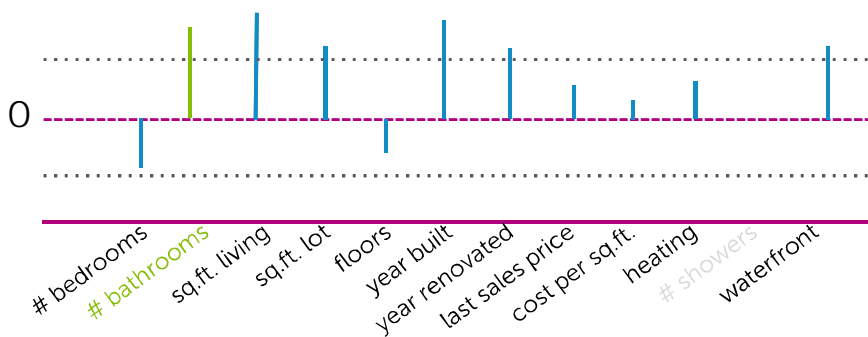
15

©2017 Emily Fox

CSF 446: Machine Learning

# Thresholding ridge coefficients?

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



16

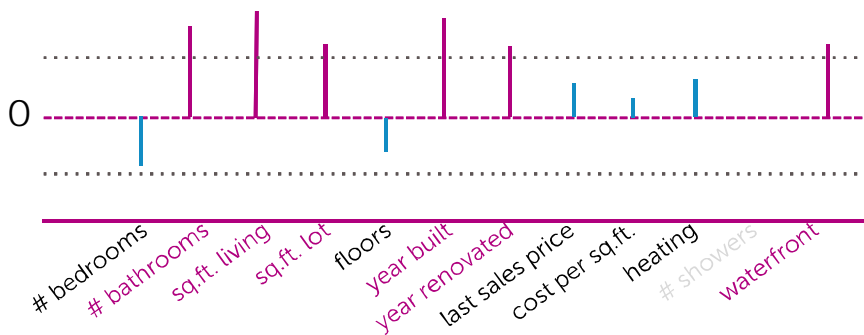
©2017 Emily Fox

CSF 446: Machine Learning



# Thresholding ridge coefficients?

Would have included bathrooms in selected model



Can regularization lead directly to sparsity?

# Try this cost instead of ridge...

Total cost =   
 $\underbrace{\text{measure of fit}}_{\text{RSS}(\mathbf{w})} + \lambda \underbrace{\text{measure of magnitude of coefficients}}_{\|\mathbf{w}\|_1 = |w_0| + \dots + |w_D|}$    
 Leads to sparse solutions!

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

## Lasso regression: $L_1$ regularized regression

Just like ridge regression, solution is governed by a continuous parameter  $\lambda$

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

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

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

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

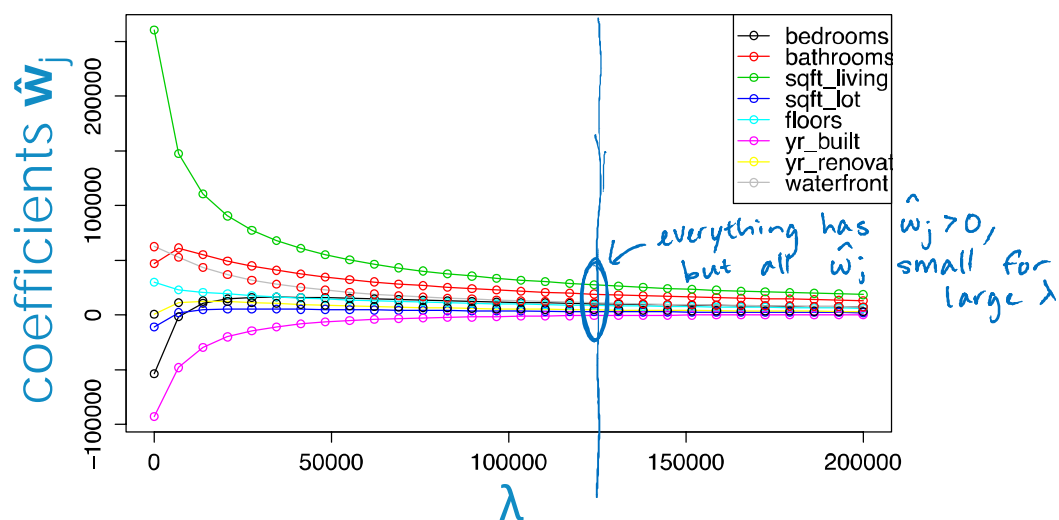
If  $\lambda$  in between:  $0 \leq \|\hat{\mathbf{w}}^{\text{lasso}}\|_1 \leq \|\hat{\mathbf{w}}^{\text{LS}}\|_1$

19

©2017 Emily Fox

CSF 446: Machine Learning

## Coefficient path – ridge

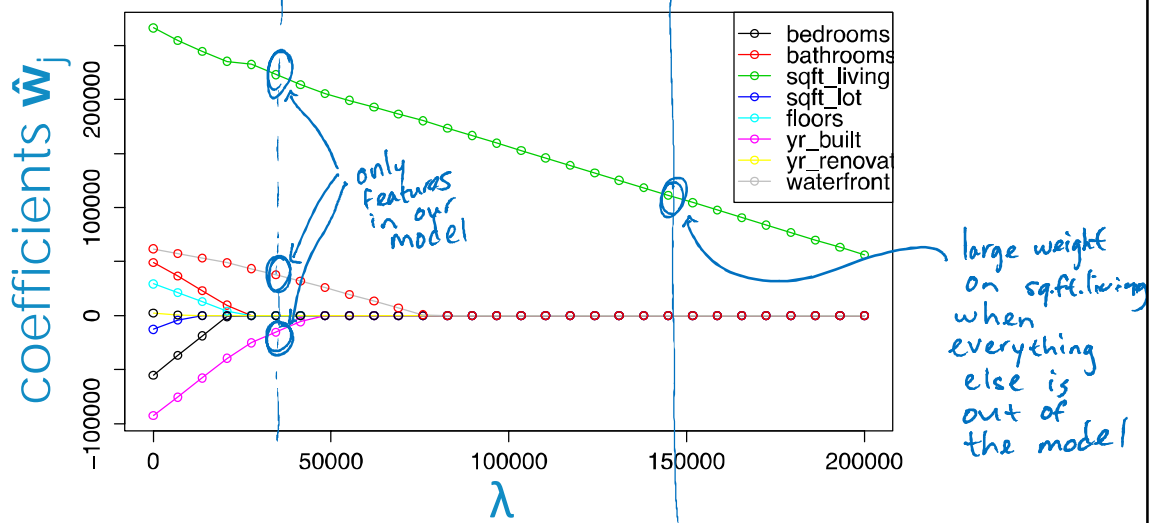


20

©2017 Emily Fox

CSF 446: Machine Learning

# Coefficient path – lasso



21

©2017 Emily Fox

CSE 446: Machine Learning

Fitting the lasso regression model  
(for given  $\lambda$  value)

## 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

23

©2017 Emily Fox

CSF 446: Machine Learning

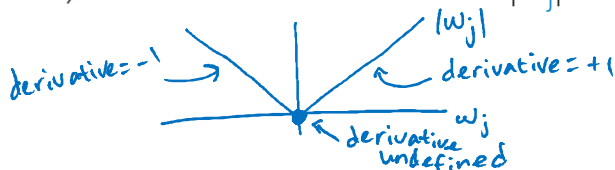
## Optimizing the lasso objective

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

$$\sum_{j=0}^D |w_j|$$

Issues:

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



gradients  $\rightarrow$  subgradients

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

can use subgradient descent

24

©2017 Emily Fox

CSF 446: Machine Learning

## Aside 1: Coordinate descent

## Coordinate descent

Goal: Minimize some function  $g$

$$\min_w g(w)$$

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

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

*when keeping others fixed*

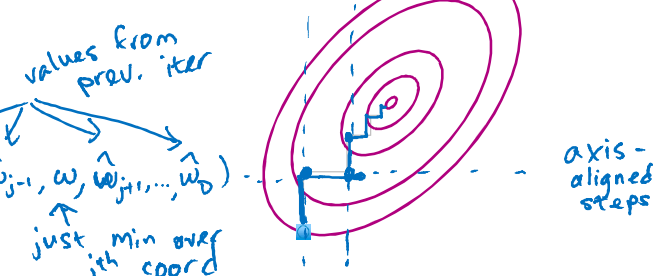
Coordinate descent:

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

while not converged

pick a coordinate  $j$

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



## 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

## Aside 2: Normalizing features

# Normalizing features

Scale training columns (not rows!) as:

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

Normalizer:  $Z_j$

Apply same training scale factors to test data:

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

Normalizer:  $Z_j$

apply to test point

summing over training points



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

## 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

1d optimization  
(coord. by coord.)

Fix all coordinates  $w_{-j}$  and take partial w.r.t.  $w_j$

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

$$= -2 \sum_{i=1}^N h_j(\mathbf{x}_i) \left( y_i - \sum_{k \neq j} w_k h_k(\mathbf{x}_i) - w_j h_j(\mathbf{x}_i) \right)$$

$$= -2 \sum_{i=1}^N h_j(\mathbf{x}_i) \left( y_i - \sum_{k \neq j} w_k h_k(\mathbf{x}_i) \right) + 2 w_j \sum_{i=1}^N h_j^2(\mathbf{x}_i)$$

$$= -2 \rho_j + 2 w_j$$

by defn of  
normalized  
feature  
= 1

31

©2017 Emily Fox

CSF 446: Machine Learning

## 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

$$\frac{\partial}{\partial w_j} \text{RSS}(\mathbf{w}) = -2 \rho_j + 2 w_j = 0$$

$$\hat{w}_j = \rho_j$$

32

©2017 Emily Fox

CSF 446: Machine Learning



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

33

©2017 Emily Fox

CSF 446: Machine Learning

## Coordinate descent for lasso (for normalized features)

## 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

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

35

©2017 Emily Fox

CSF 446: Machine Learning

## Coordinate descent for lasso

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

while not converged

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

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

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

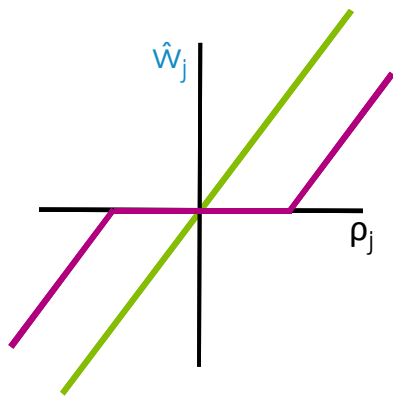
36

©2017 Emily Fox

CSF 446: Machine Learning

## Soft thresholding

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



37

©2017 Emily Fox

CSF 446: Machine Learning

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

38

©2017 Emily Fox

CSF 446: Machine Learning

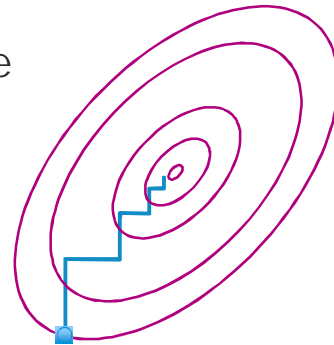
## Convergence criteria

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  $\text{max step} < \epsilon$



39

©2017 Emily Fox

CSF 446: Machine Learning

## 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]

40

©2017 Emily Fox

CSF 446: Machine Learning

## Coordinate descent for lasso (for unnormalized features)

## Coordinate descent for lasso with normalized features

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

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

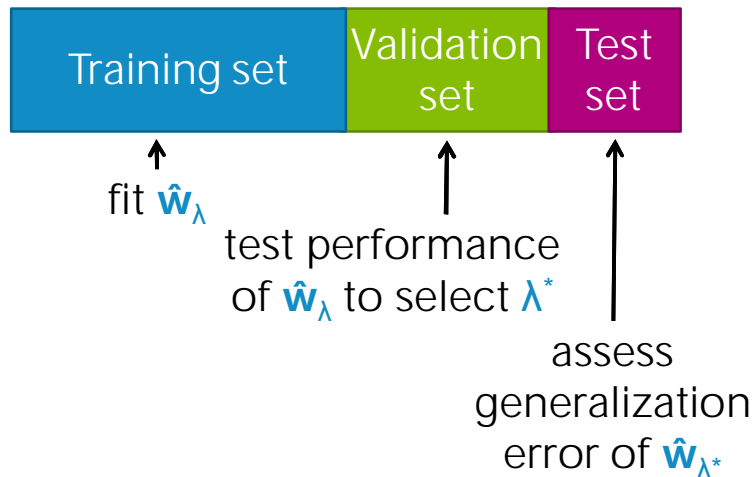
43

©2017 Emily Fox

CSF 446: Machine Learning

How to choose  $\lambda$

If sufficient amount of data...



45

©2017 Emily Fox

CSF 446: Machine Learning

Summary for feature selection  
and lasso regression

46

## 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

47

©2017 Emily Fox

CSF 446: 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  $\lambda$  is varied
- Interpret lasso coefficient path plot
- Contrast ridge and lasso regression
- Estimate lasso regression parameters using an iterative coordinate descent algorithm

48

©2017 Emily Fox

CSF 446: Machine Learning



## Deriving the lasso coordinate descent update

**OPTIONAL**



## 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  $w_{-j}$  and take partial w.r.t.  $w_j$

derive without normalizing features

## Part 1: Partial of RSS term

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

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

51

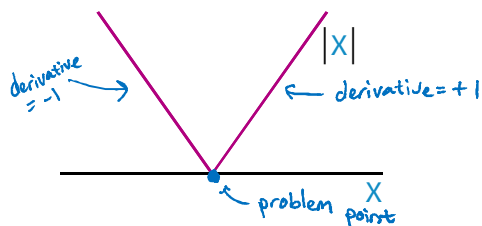
©2017 Emily Fox

CSF 446: Machine Learning

## Part 2: Partial of $L_1$ penalty term

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

$$\lambda \frac{\partial}{\partial w_j} |w_j| = ???$$



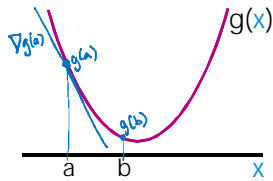
52

©2017 Emily Fox

CSF 446: Machine Learning

# Subgradients of convex functions

Gradients lower bound convex functions:

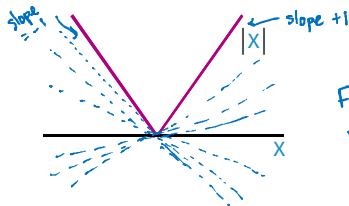


$$g(b) \geq g(a) + \nabla g(a)(b-a)$$

unique at x if function differentiable at x

Subgradients: Generalize gradients to non-differentiable points:

- Any plane that lower bounds function



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

$$V \in \partial g(x) \text{ subgradient of } g \text{ at } x \text{ if } g(b) \geq g(a) + V(b-a)$$

53

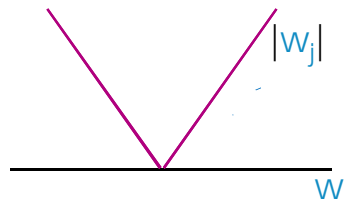
©2017 Emily Fox

CSF 446: Machine Learning

## Part 2: Subgradient of $L_1$ term

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

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



54

©2017 Emily Fox

CSF 446: Machine Learning

## Putting it all together...

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

from RSS
from  $\lambda \|\cdot\|_1$  penalty

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

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

55

©2017 Emily Fox

CSF 446: Machine Learning

## 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 ( $w_j < 0$ ):  $2z_j \hat{w}_j - 2\rho_j - \lambda = 0$   
 $\hat{w}_j = \frac{2\rho_j + \lambda}{2z_j} = \frac{\rho_j + \frac{\lambda}{2}}{z_j}$

For  $\hat{w}_j < 0$ , need  
 $\rho_j < -\frac{\lambda}{2}$

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:  
 $-2\rho_j + \lambda \geq 0 \rightarrow \rho_j \leq \frac{\lambda}{2}$      $-2\rho_j - \lambda \leq 0 \rightarrow \rho_j \geq -\frac{\lambda}{2}$      $-\frac{\lambda}{2} \leq \rho_j \leq \frac{\lambda}{2}$   
so that  $\hat{w}_j = 0$  is an optimum

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

For  $\hat{w}_j > 0$ , need  
 $\rho_j > \frac{\lambda}{2}$

56

©2017 Emily Fox

CSF 446: Machine Learning

Optimal solution:  
Set subgradient = 0

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

$= 0$



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

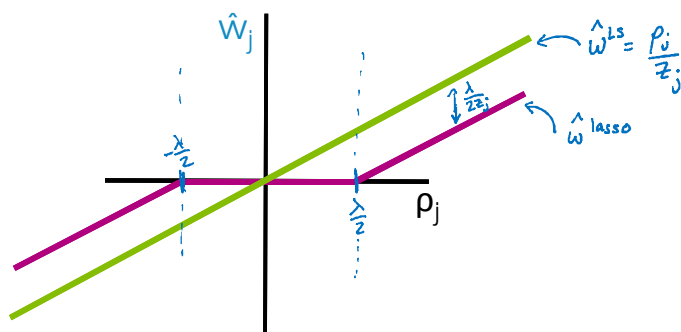
57

©2017 Emily Fox

CSF 446: Machine Learning

## Soft thresholding

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



58

©2017 Emily Fox

CSF 446: Machine Learning

## Coordinate descent for lasso

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

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