

# Simple Variable Selection LASSO: Sparse Regression

Machine Learning – CSE546  
 Carlos Guestrin  
 University of Washington  
 October 7, 2014

©2005-2014 Carlos Guestrin

1

## Sparsity

*100M - 100Bs of params*

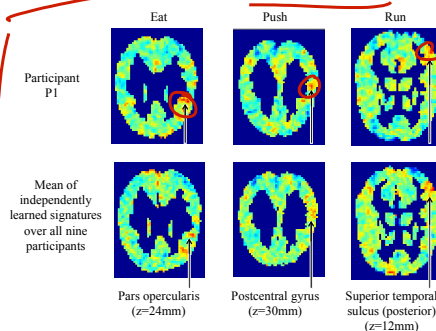
- Vector  $\mathbf{w}$  is sparse, if many entries are zero:
  - ↳ {1.7, 2.2, 0.0, 0.0, 2.9, 0.0, 0.0, ..., -3.1, ...}*
- Very useful for many tasks, e.g.,
  - Efficiency:** If  $\text{size}(\mathbf{w}) = 100B$ , each prediction is expensive:
    - If part of an online system, too slow
    - If  $\mathbf{w}$  is sparse, prediction computation only depends on number of non-zeros

*Best are zero*  
 $f(x) = w_0 + \sum_{i=1}^k w_i h_i(x)$

- Interpretability:** What are the relevant dimension to make a prediction?
  - E.g., what are the parts of the brain associated with particular words?

*with K non-zeros*  
*(100B / K) Subsets*

- But computationally intractable to perform “all subsets” regression



©2005-2014 Carlos Guestrin

2

## Simple greedy model selection algorithm

- Pick a dictionary of features
  - e.g., polynomials for linear regression
- Greedy heuristic:
  - Start from empty (or simple) set of features  $F_0 = \emptyset$ , or the constant  $w_0$
  - Run learning algorithm for current set of features  $F_t$ 
    - Obtain  $h_t$
  - Select **next best feature**  $X_i^*$ 
    - e.g.,  $X_i$  that results in lowest training error learner when learning with  $F_t + \{X_i\}$
  - $F_{t+1} \leftarrow F_t + \{X_i^*\}$
  - Recurse

runs of learning alg. !!  
1008

©2005-2014 Carlos Guestrin

3

## Greedy model selection

- Applicable in many settings:
  - Linear regression: Selecting basis functions
  - Naïve Bayes: Selecting (independent) features  $P(X_i|Y)$
  - Logistic regression: Selecting features (basis functions)
  - Decision trees: Selecting leaves to expand
- Only a heuristic!
  - But, sometimes you can prove something cool about it
    - e.g., [Krause & Guestrin '05]: Near-optimal in some settings that include Naïve Bayes
- There are many more elaborate methods out there

©2005-2014 Carlos Guestrin

4

# When do we stop???

- Greedy heuristic:

- ...
- Select **next best feature**  $X_i^*$ 
  - e.g.,  $X_j$  that results in lowest training error learner when learning with  $F_t + \{X_j\}$
- $F_{t+1} \leftarrow F_t + \{X_i^*\}$
- Recurse

*either k features*

**When do you stop???**

- ~~When training error is low enough?~~
- ~~When test set error is low enough?~~
- *cross validation please!!*

# Regularization in Linear Regression

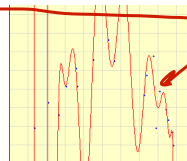
- Overfitting usually leads to very large parameter choices, e.g.:

$-2.2 + 3.1 X - 0.30 X^2$



*simple models*

$-1.1 + 4,700,910.7 X - 8,585,638.4 X^2 + \dots$



*overfitting*

*penalty for large weights*

- Regularized** or **penalized** regression aims to impose a “complexity” penalty by penalizing large weights

- “Shrinkage” method

$L_2$  regularization

*penalizes towards smoother functions*

# Variable Selection by Regularization

- Ridge regression: Penalizes large weights *with L2 norm → smooth functions*
- What if we want to perform “feature selection”?
  - E.g., Which regions of the brain are important for word prediction?
  - Can't simply choose features with largest coefficients in ridge solution  
*lots of small, but non-zero coeffs → Ridge  
 (0.001, 0.1, 3.2, -2.7, -0.02)*
- Try new penalty: Penalize non-zero weights
  - Regularization penalty:  $\|w\|_1 = \sum_i |w_i|$   
*LASSO*
  - Leads to sparse solutions
  - Just like ridge regression, solution is indexed by a continuous param  $\lambda$
  - This simple approach has changed statistics, machine learning & electrical engineering

©2005-2014 Carlos Guestrin

7

# LASSO Regression

- **LASSO**: least absolute shrinkage and selection operator

- New objective:

$$\min_w \sum_{j=1}^N \left( f(x_j) - \left( w_0 + \sum_{i=1}^k w_i h_i(x_j) \right) \right)^2 + \lambda \sum_{i=1}^k |w_i|$$

*RSS*

*truth*      *f*

↑  
 please don't regularize  $w_0$ , it didn't do anything to you

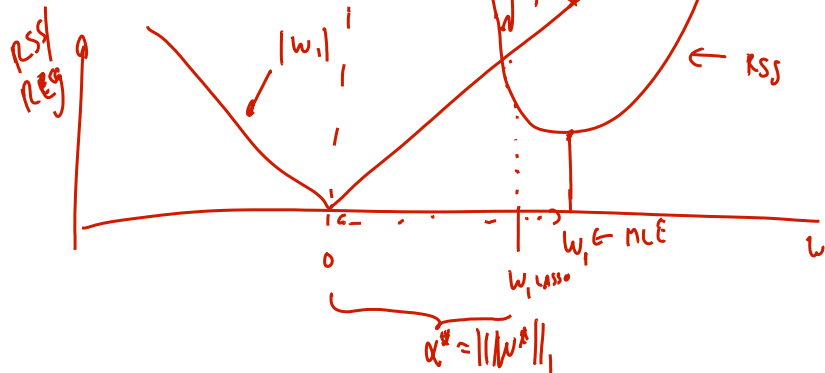
©2005-2014 Carlos Guestrin

8

# Geometric intuition of regularized objectives in 1d

- LASSO solution:

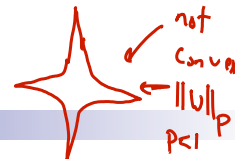
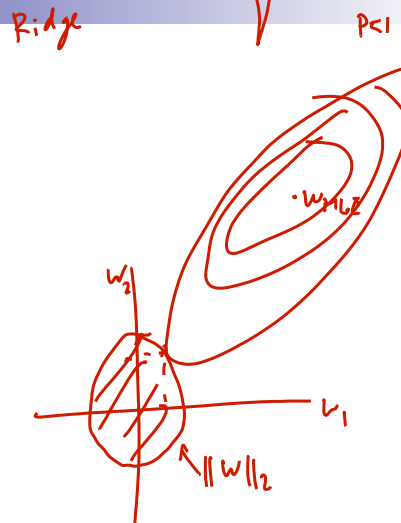
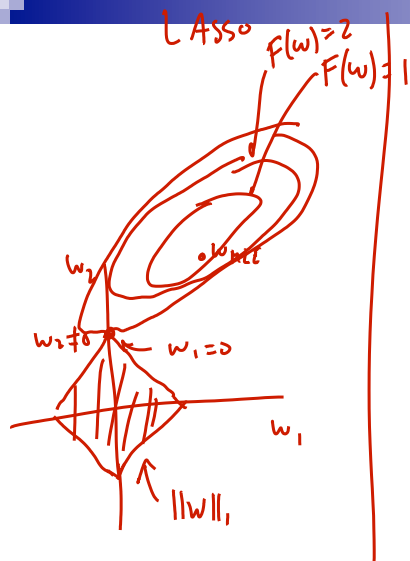
$$\hat{w}_{LASSO} = \arg \min_w \sum_{j=1}^N \left( t(x_j) - (w_0 + \sum_{i=1}^k w_i h_i(x_j)) \right)^2 + \lambda \sum_{i=1}^k |w_i|$$



©2005-2014 Carlos Guestrin

9

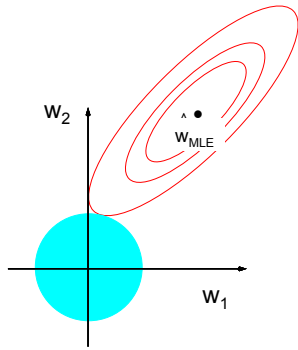
## Geometric intuition



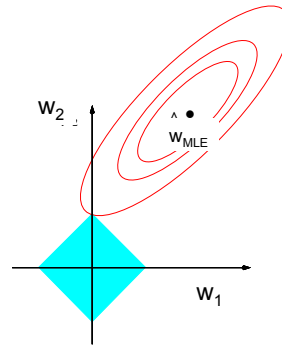
©2005-2014 Carlos Guestrin

10

# Geometric Intuition for Sparsity



Ridge Regression



Lasso

From Rob Tibshirani slides

# Optimizing the LASSO Objective



- LASSO solution:

$$\hat{w}_{LASSO} = \arg \min_w \sum_{j=1}^N \left( t(x_j) - (w_0 + \sum_{i=1}^k w_i h_i(x_j)) \right)^2 + \lambda \sum_{i=1}^k |w_i|$$

Take the derivative & set to 0

1. Derivative of  $|w_i|$ ?



2. Even if you could take derivatives, no closed-form solution exists.

# Coordinate Descent



- Given a function  $F$ 
  - Want to find minimum  $\hat{w} = \text{argmin} F(w_0, w_1, \dots, w_k)$
- Often, hard to find minimum for all coordinates, but easy for one coordinate
- Coordinate descent:
  - initialize  $w=0$  or something smarter
  - while not converged
    - Pick a coordinate  $l$
    - $\hat{w}_l = \text{argmin}_{w_l} F(\hat{w}_0, \hat{w}_1, \dots, \hat{w}_{l-1}, w_l, \hat{w}_{l+1}, \dots, \hat{w}_k)$ 
      - Fix to values from previous iteration
- How do we pick next coordinate?
  - random, round robin, "smartly"
  - converges!!
  - but usually to a local optima
- Super useful approach for \*many\* problems
  - Converges to optimum in some cases, such as LASSO

©2005-2014 Carlos Guestrin

13

# Optimizing LASSO Objective One Coordinate at a Time

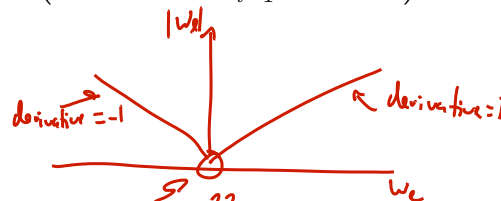
$$\sum_{j=1}^N \left( t(x_j) - \left( w_0 + \sum_{i=1}^k w_i h_i(x_j) \right) \right)^2 + \lambda \sum_{i=1}^k |w_i|$$

- Taking the derivative:
  - Residual sum of squares (RSS):

$$\frac{\partial}{\partial w_l} \text{RSS}(\mathbf{w}) = -2 \sum_{j=1}^N h_l(x_j) \left( t(x_j) - \left( w_0 + \sum_{i=1}^k w_i h_i(x_j) \right) \right)$$

- Penalty term:

$$\frac{\partial}{\partial w_k} \lambda \sum_{i=1}^k |w_i| = \lambda \frac{\partial}{\partial w_k} |w_k|$$



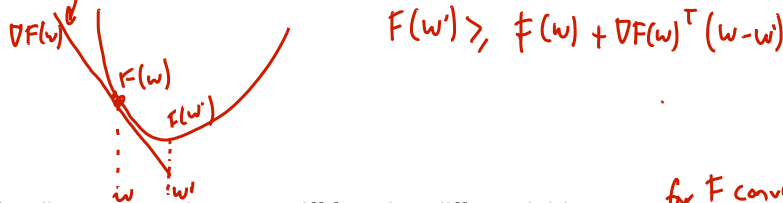
at opt solution  
most params are 0

©2005-2014 Carlos Guestrin

14

# Subgradients of Convex Functions

- Gradients lower bound convex functions:

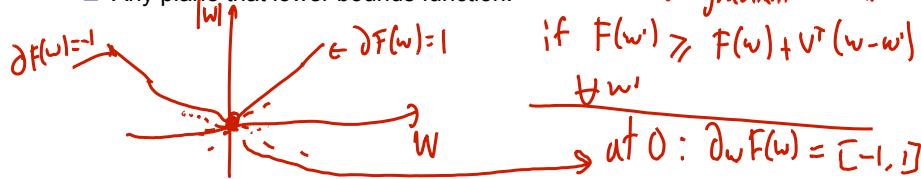


- Gradients are unique at  $w$  iff function differentiable at  $w$

for  $F$  convex

- Subgradients: Generalize gradients to non-differentiable points:

- Any plane that lower bounds function:



$v$  is a subgradient of  $F$  at  $w$  if  $F(w') \geq F(w) + v^T (w-w')$

at  $0: \partial_w F(w) = [-1, 1]$

©2005-2014 Carlos Guestrin

15

At optimum, subgradient = 0

## Taking the Subgradient

$$\sum_{j=1}^N \left( t(x_j) - (w_0 + \sum_{i=1}^k w_i h_i(x_j)) \right)^2 + \lambda \sum_{i=1}^k |w_i|$$

- Gradient of RSS term:

$$a_\ell = 2 \sum_{j=1}^N (h_\ell(x_j))^2 > 0$$

$$\frac{\partial}{\partial w_\ell} RSS(\mathbf{w}) = a_\ell w_\ell - c_\ell$$

$$c_\ell = 2 \sum_{j=1}^N h_\ell(x_j) \left( t(x_j) - (w_0 + \sum_{i \neq \ell} w_i h_i(x_j)) \right)$$

- If no penalty:

$$w_\ell = c_\ell / a_\ell$$

- Subgradient of full objective:

$$\partial_{w_\ell} F(\mathbf{w}) = a_\ell w_\ell - c_\ell + \lambda \partial_{w_\ell} |w_\ell|$$

$$\begin{cases} -1 & \text{when } w_\ell < 0 \\ [-1, 1] & \text{when } w_\ell = 0 \\ +1 & \text{when } w_\ell > 0 \end{cases}$$

$$= \begin{cases} a_\ell w_\ell - c_\ell - \lambda & \text{when } w_\ell < 0 \\ [-c_\ell - \lambda, -c_\ell + \lambda] & \text{when } w_\ell = 0 \\ a_\ell w_\ell - c_\ell + \lambda & \text{when } w_\ell > 0 \end{cases} \quad \left. \begin{array}{l} \text{at the OPT we} \\ \text{this beast is } \emptyset \end{array} \right\}$$

©2005-2014 Carlos Guestrin

16



# Setting Subgradient to 0

$a_l > 0$

Choose  $w_l$  such that

$$0 \in \partial_{w_l} F(\mathbf{w}) = \begin{cases} a_l w_l - c_l - \lambda & w_l < 0 \\ [-c_l - \lambda, -c_l + \lambda] & w_l = 0 \\ a_l w_l - c_l + \lambda & w_l > 0 \end{cases}$$

Optimum with  $w_l < 0$  ?  $\Rightarrow a_l w_l - c_l - \lambda = 0 \Rightarrow w_l = \frac{c_l + \lambda}{a_l}$  when  $c_l + \lambda < 0$

Optimum with  $w_l > 0$  ?  $\Rightarrow a_l w_l - c_l + \lambda = 0 \Rightarrow w_l = \frac{c_l - \lambda}{a_l}$  when  $c_l - \lambda > 0$

Optimum with  $w_l = 0$   $\Rightarrow 0 \in [-c_l - \lambda, -c_l + \lambda] \Rightarrow -\lambda \leq c_l \leq \lambda$

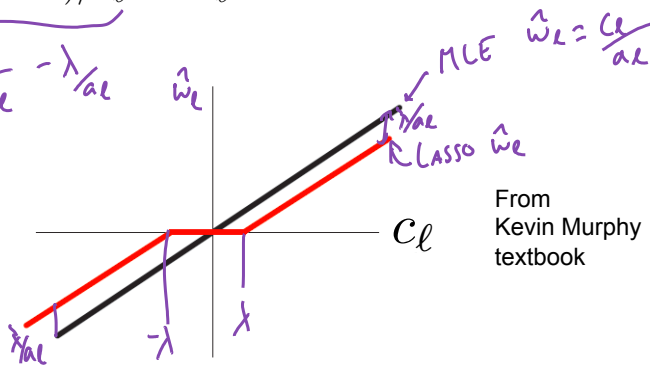
this is why we get sparsity  $\Rightarrow w_l = 0$   
 ↳ sparsity in coord  $l$

# Soft Thresholding

Reminder, for MLE

$$\hat{w}_l = \frac{c_l}{a_l}$$

$$\hat{w}_l = \begin{cases} (c_l + \lambda)/a_l & c_l < -\lambda \\ 0 & c_l \in [-\lambda, \lambda] \\ (c_l - \lambda)/a_l & c_l > \lambda \end{cases}$$



From Kevin Murphy textbook

# Coordinate Descent for LASSO (aka Shooting Algorithm)

- Repeat until convergence

- Pick a coordinate  $\ell$  at (random or sequentially)

minimum F  
w.r.t  $w_\ell$   
using subgradient

$$\hat{w}_\ell = \begin{cases} (c_\ell + \lambda)/a_\ell & c_\ell < -\lambda \\ 0 & c_\ell \in [-\lambda, \lambda] \\ (c_\ell - \lambda)/a_\ell & c_\ell > \lambda \end{cases}$$

- Where:

$$a_\ell = 2 \sum_{j=1}^N (h_\ell(x_j))^2$$

$$c_\ell = 2 \sum_{j=1}^N h_\ell(x_j) \left( t(x_j) - (\hat{w}_0 + \sum_{i \neq \ell} \hat{w}_i h_i(x_j)) \right)$$

value in previous iteration  
of all  $\hat{w}_i$  except for  $\ell$

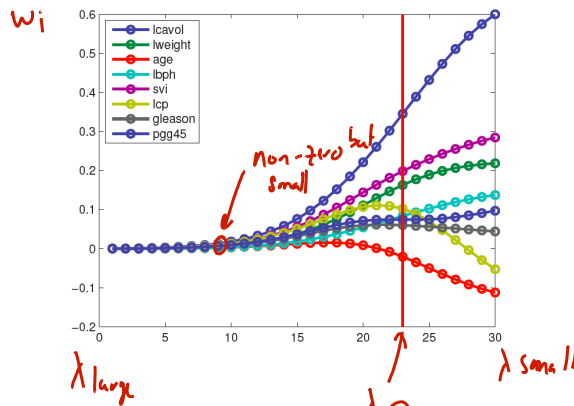
- For convergence rates, see Shalev-Shwartz and Tewari 2009

- Other common technique = LARS

- Least angle regression and shrinkage, Efron et al. 2004

what about  
 $w_0$ ?  
don't regularize  
 $\hat{w}_0 = c_0 / a_0$

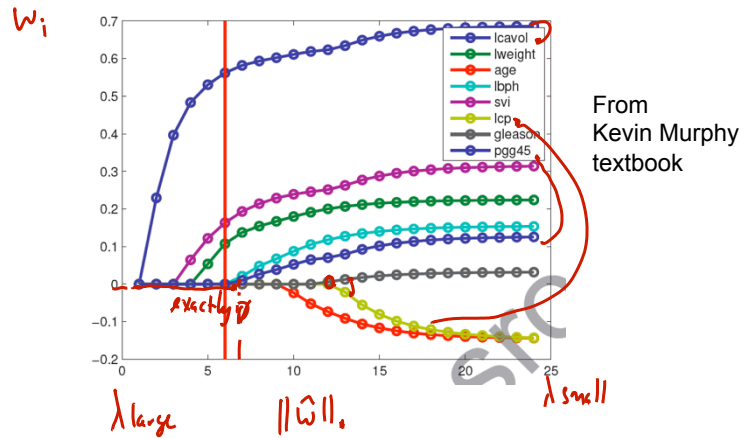
# Recall: Ridge Coefficient Path



$\|w\|_2^2$  reg.

- Typical approach: select  $\lambda$  using cross validation

## Now: LASSO Coefficient Path



©2005-2014 Carlos Guestrin

21

## LASSO Example

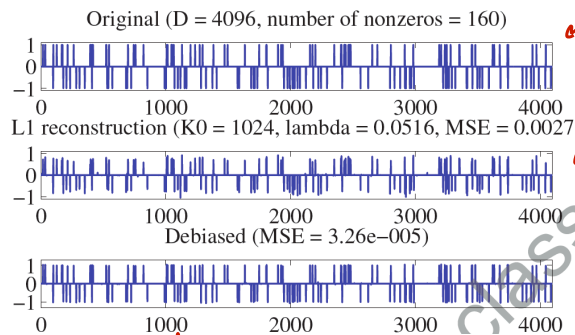
Term	Least Squares	Ridge	Lasso
Intercept	2.465	2.452	2.468
lcavol	0.680	0.420	0.533
lweight	0.263	0.238	0.169
age	-0.141	-0.046	
lbph	0.210	0.162	0.002
svi	0.305	0.227	0.094
lcp	-0.288	0.000	
gleason	-0.021	0.040	
pgg45	0.267	0.133	

From Rob Tibshirani slides

©2005-2014 Carlos Guestrin

22

# Debiasing



From Kevin Murphy textbook

©2005-2014 Carlos Guestrin

23

# What you need to know

- Variable Selection: find a sparse solution to learning problem
- $L_1$  regularization is one way to do variable selection
  - Applies beyond regressions
  - Hundreds of other approaches out there
- LASSO objective non-differentiable, but convex → Use subgradient
- No closed-form solution for minimization → Use coordinate descent
- Shooting algorithm is very simple approach for solving LASSO

©2005-2014 Carlos Guestrin

24