

The regression problem in matrix notation
$$\mathbf{w}^* = \arg\min_{\mathbf{w}} \sum_{j \in I} \left(t(\mathbf{x}_j) - \sum_{i \in I} w_i h_i(\mathbf{x}_j) \right)^2$$

$$\mathbf{w}^* = \arg\min_{\mathbf{w}} \left(\mathbf{H}\mathbf{w} - \mathbf{t} \right)^T (\mathbf{H}\mathbf{w} - \mathbf{t})$$

$$\operatorname{residual error}$$

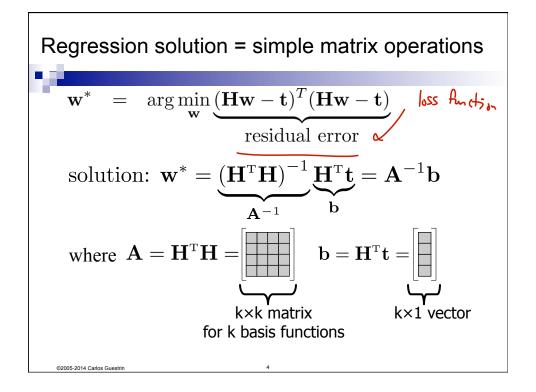
$$\mathbf{w}$$

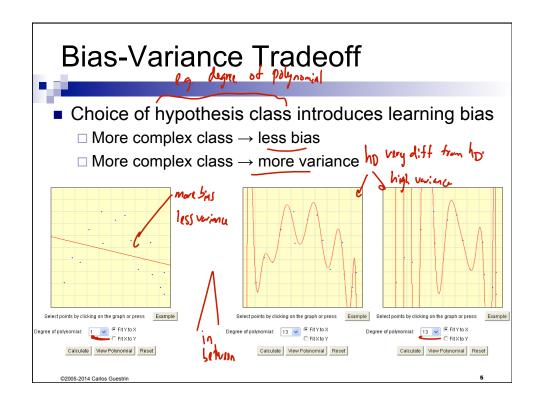
$$\mathbf{w} = \lim_{\mathbf{x} \in I} \left(\mathbf{x}_j \right) - \sum_{i \in I} w_i h_i(\mathbf{x}_j) \right)$$

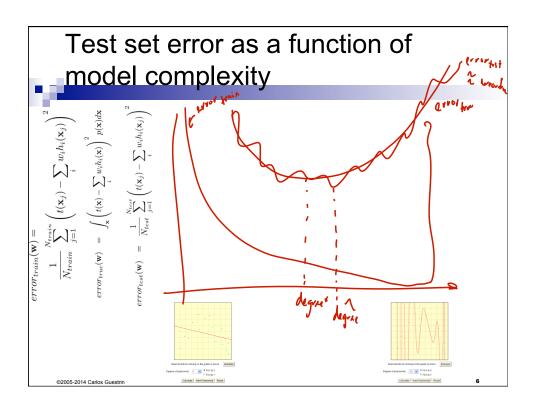
$$\operatorname{residual error}$$

$$\operatorname{residual error}$$

$$\mathbf{w} = \lim_{\mathbf{x} \in I} \left(\mathbf{x}_j \right) - \lim_{\mathbf{x} \in I} \left$$



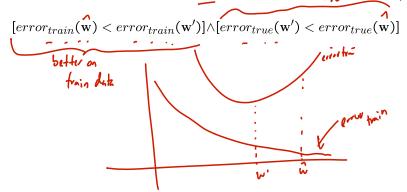




Overfitting



■ Overfitting: a learning algorithm overfits the training data if it outputs a solution w when there exists another solution w such that:



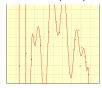
©2005-2014 Carlos Guestria

Regularization in Linear Regression



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

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



- Regularized or penalized regression aims to impose a "complexity" penalty by penalizing large weights
 - □ "Shrinkage" method

©2005-2013 Carlos Guestrin

Quadratic Penalty (regularization)



- What we thought we wanted to minimize:
- But weights got too big, penalize large weights:

©2005-2013 Carlos Guestrir

.

Ridge Regression



- Ameliorating issues with overfitting:
- New objective:

©2005-2013 Carlos Guestrin

Ridge Regression in Matrix Notation

$$\hat{\mathbf{w}}_{ridge} = \arg\min_{w} \sum_{i=1}^{N} \left(t(x_i) - (w_0 + \sum_{i=1}^{k} w_i h_i(x_j)) \right)^2 + \lambda \sum_{i=1}^{k} w_i^2$$

$$= \underset{\mathbf{w}}{\operatorname{arg\,min}} \underbrace{(\mathbf{H}\mathbf{w} - \mathbf{t})^{T}(\mathbf{H}\mathbf{w} - \mathbf{t})}_{\text{residual error}} + \lambda \ \mathbf{w}^{T} I_{0+k} \mathbf{w}$$

Minimizing the Ridge Regression Objective

$$\mathbf{w}_{\text{\tiny MLE}}^* = \underbrace{\left(\mathbf{H}^{\text{\tiny T}}\mathbf{H}\right)^{-1}}_{\mathbf{A}^{-1}} \underbrace{\mathbf{H}^{\text{\tiny T}}\mathbf{t}}_{\mathbf{b}} = \mathbf{A}^{-1}\mathbf{b}$$

$$\mathbf{w}_{\text{MLE}}^* = \underbrace{\left(\mathbf{H}^{\text{T}}\mathbf{H}\right)^{-1}}_{\mathbf{A}^{-1}} \underbrace{\mathbf{H}^{\text{T}}\mathbf{t}}_{\mathbf{b}} = \mathbf{A}^{-1}\mathbf{b}$$

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

$$= (H\mathbf{w} - \mathbf{t})^T (H\mathbf{w} - \mathbf{t}) + \lambda \mathbf{w}^T I_{0+k} \mathbf{w}$$

Shrinkage Properties



$$\hat{\mathbf{w}}_{ridge} = (H^T H + \lambda \ I_{0+k})^{-1} H^T \mathbf{t}$$

lacksquare If orthonormal features/basis: $H^T H = I$

©2005-2013 Carlos Guestri

40

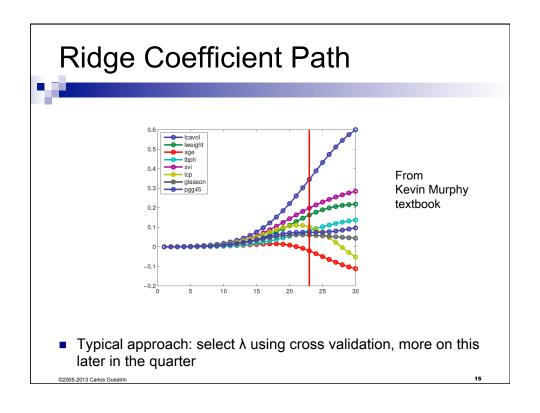
Ridge Regression: Effect of Regularization

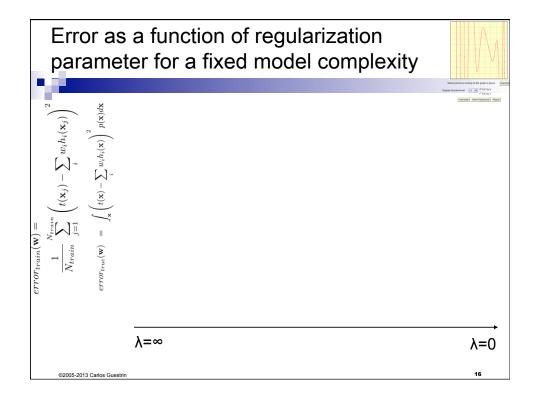


$$\hat{\mathbf{w}}_{ridge} = \arg\min_{w} \sum_{j=1}^{N} \left(t(x_j) - (w_0 + \sum_{i=1}^{k} w_i h_i(x_j)) \right)^2 + \lambda \sum_{i=1}^{k} w_i^2$$

- Solution is indexed by the regularization parameter λ
- Larger λ
- Smaller λ
- As $\lambda \rightarrow 0$
- As λ →∞

©2005-2013 Carlos Guestrin

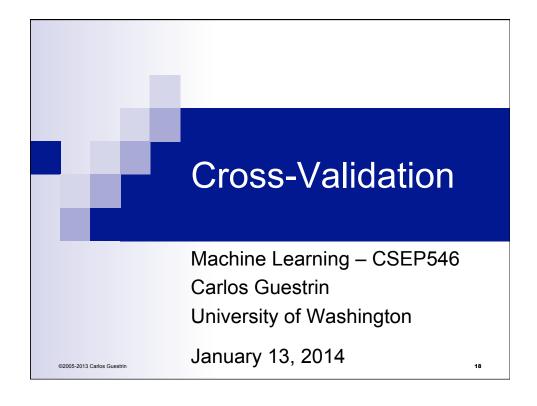


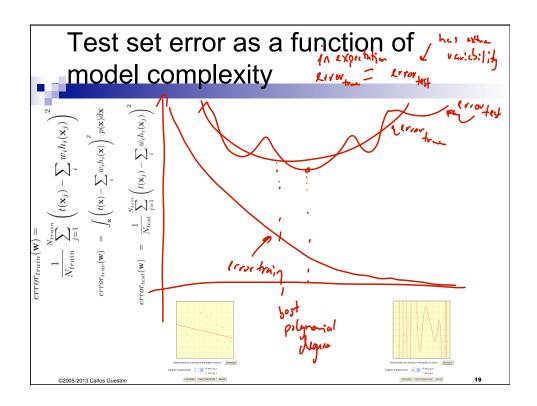


What you need to know...

- - Regularization
 - □ Penalizes for complex models
 - Ridge regression
 - □ L₂ penalized least-squares regression
 - □ Regularization parameter trades off model complexity with training error

©2005-2013 Carlos Guestrin





How... How???????

- - How do we pick the regularization constant λ...
 - □ And all other constants in ML, 'cause one thing ML doesn't lack is constants to tune… ⑤
 - We could use the test data, but...

©2005-2013 Carlos Guestrin

(LOO) Leave-one-out cross validation



- Consider a validation set with 1 example:
 - □ D training data
 - \Box D\j training data with jth data point moved to validation set
- Learn classifier $h_{D\setminus j}$ with $D\setminus j$ dataset
- Estimate true error as squared error on predicting t(x_i):
 - □ Unbiased estimate of $error_{true}(\boldsymbol{h}_{\boldsymbol{D}\setminus i})!$
 - □ Seems really bad estimator, but wait!
- LOO cross validation: Average over all data points *j*:
 - $\ \square$ For each data point you leave out, learn a new classifier $h_{D_{||}}$
 - Estimate error as: $error_{LOO} = \frac{1}{N} \sum_{j=1}^{N} \left(t(\mathbf{x}_j) h_{\mathcal{D} \backslash j}(\mathbf{x}_j) \right)^2$

©2005-2013 Carlos Guestrir

21

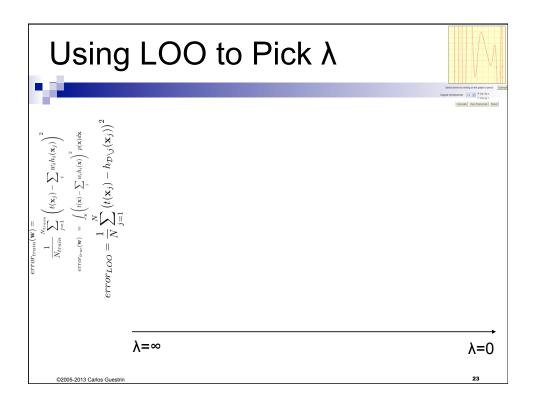
LOO cross validation is (almost) unbiased estimate of true error of h_D !

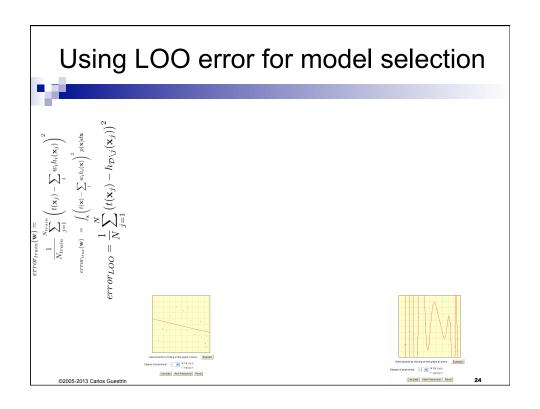


- When computing LOOCV error, we only use N-1 data points
 - □ So it's not estimate of true error of learning with *N* data points!
 - □ Usually pessimistic, though learning with less data typically gives worse answer
- LOO is almost unbiased!

- Great news!
 - ☐ Use LOO error for model selection!!!
 - E.g., picking λ

©2005-2013 Carlos Guestrin





Computational cost of LOO



- Suppose you have 100,000 data points
- You implemented a great version of your learning algorithm
 - □ Learns in only 1 second
- Computing LOO will take about 1 day!!!
 - ☐ If you have to do for each choice of basis functions, it will take fooooooreeeve'!!!
- Solution 1: Preferred, but not usually possible
 - ☐ Find a cool trick to compute LOO (e.g., see homework)

Solution 2 to complexity of computing LOO:

(More typical) Use k-fold cross validation

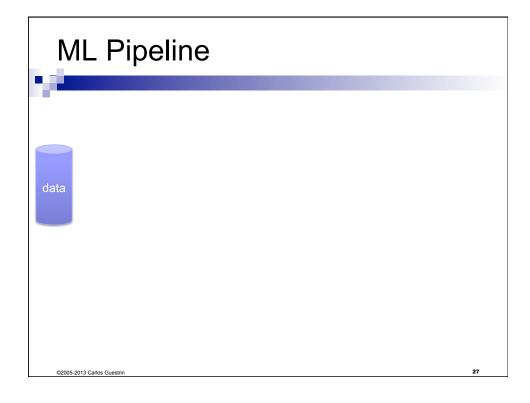


- Randomly divide training data into k equal parts
 - $\square D_1,...,D_k$
- For each i
 - □ Learn classifier $h_{D \setminus D_i}$ using data point not in D_i

• k-fold cross validation error is average over data splits:

$$error_{k-fold} = \frac{1}{k} \sum_{i=1}^{k} error_{\mathcal{D}_i}$$

- k-fold cross validation properties:
 - Much faster to compute than LOO
 - □ More (pessimistically) biased using much less data, only m(k-1)/k
 - □ Usually, k = 10 ②

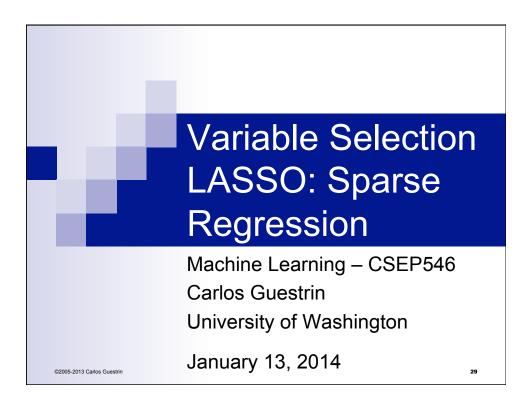


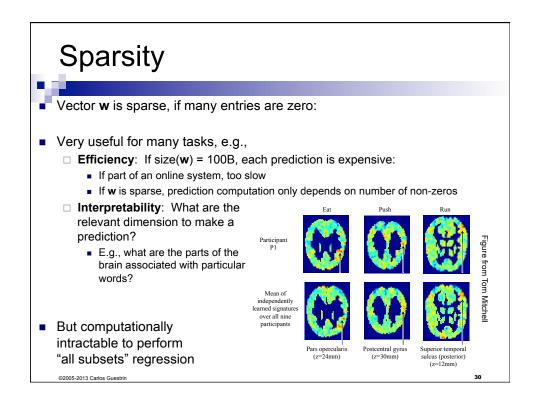
What you need to know...



- Use cross-validation to choose magic parameters such as λ
- Leave-one-out is the best you can do, but sometimes too slow
 - ☐ In that case, use k-fold cross-validation

©2005-2013 Carlos Guestrin





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 F₀ = Ø
 - □ Run learning algorithm for current set of features F_t
 - Obtain *h*,
 - ☐ Select next best feature X_i*
 - e.g., X_j that results in lowest training error learner when learning with F_t + {X_i}
 - $\Box F_{t+1} \leftarrow F_t + \{X_i^*\}$
 - □ Recurse

©2005-2013 Carlos Guestri

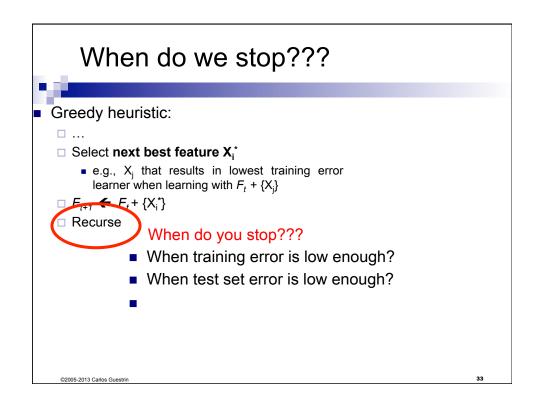
31

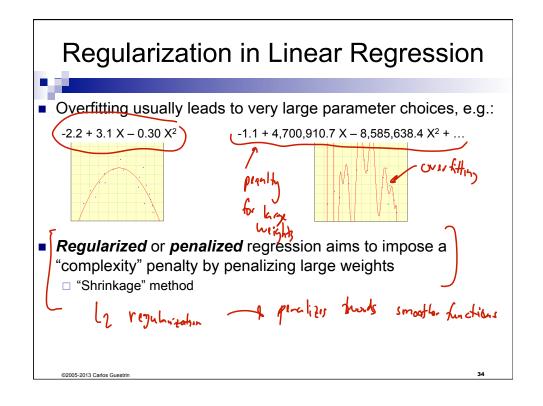
Greedy model selection



- 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

©2005-2013 Carlos Guestria





Variable Selection by Regularization



- Ridge regression: Penalizes large weights
- 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
- Try new penalty: Penalize non-zero weights
 - □ Regularization penalty:
 - □ Leads to sparse solutions
 - $\hfill \square$ Just like ridge regression, solution is indexed by a continuous param λ
 - ☐ This simple approach has changed statistics, machine learning & electrical engineering

©2005-2013 Carlos Guestrin

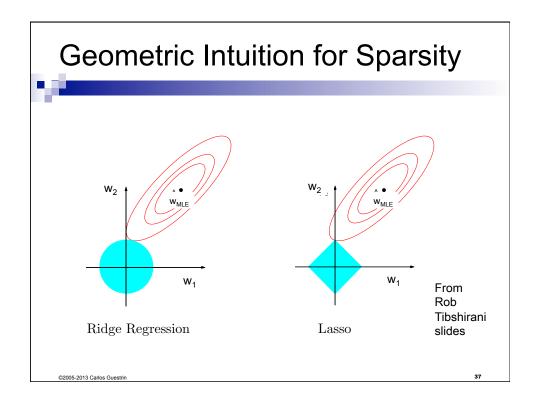
35

LASSO Regression



- LASSO: least absolute shrinkage and selection operator
- New objective:

©2005-2013 Carlos Guestria



Optimizing the LASSO Objective

LASSO solution:

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

©2005-2013 Carlos Guestrin

Coordinate Descent

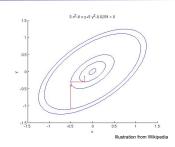
- N
- Given a function F
 - □ Want to find minimum
- Often, hard to find minimum for all coordinates, but easy for one coordinate
- Coordinate descent:
- How do we pick next coordinate?
- Super useful approach for *many* problems
 - $\hfill\Box$ Converges to optimum in some cases, such as LASSO

39

How do we find the minimum over each coordinate?



- Key step in coordinate descent:
 - □ Find minimum over each coordinate



Standard approach:

©2005-2013 Carlos Guestrin

Optimizing LASSO Objective One Coordinate at a Time



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

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

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

□ Penalty term:

Coordinate Descent for LASSO (aka Shooting Algorithm)



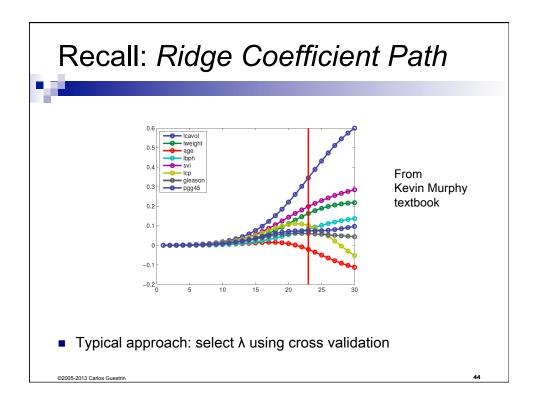
- Repeat until convergence
 - □ Pick a coordinate *l* at (random or sequentially)
 - et: $\hat{w}_{\ell} = \left\{ \begin{array}{ll} (c_{\ell} + \lambda)/a_{\ell} & c_{\ell} < -\lambda \\ 0 & c_{\ell} \in [-\lambda, \lambda] \\ (c_{\ell} \lambda)/a_{\ell} & c_{\ell} > \lambda \end{array} \right.$

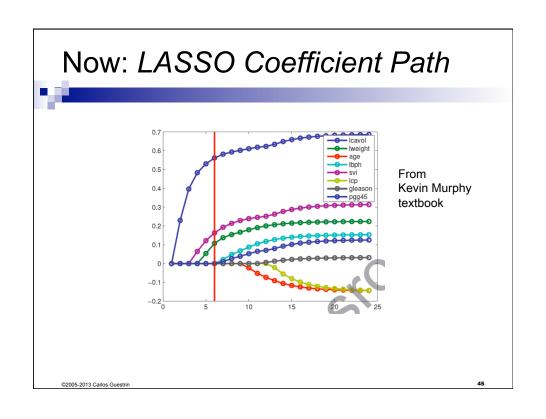
■ Where:
$$a_{\ell} = 2 \sum_{j=1}^{N} (h_{\ell}(\mathbf{x}_{j}))^{2}$$

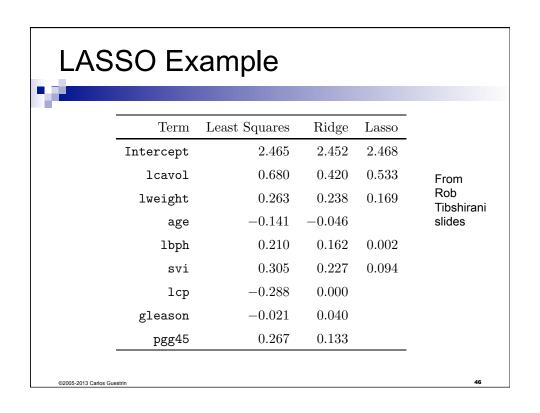
$$c_{\ell} = 2 \sum_{j=1}^{N} h_{\ell}(\mathbf{x}_{j}) \left(t(\mathbf{x}_{j}) - (w_{0} + \sum_{i \neq \ell} w_{i} h_{i}(\mathbf{x}_{j})) \right)$$

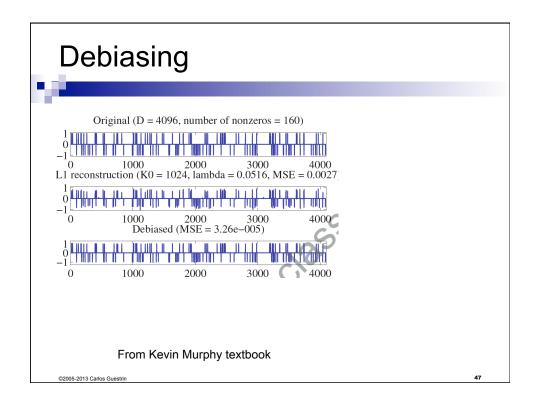
- ☐ For convergence rates, see Shalev-Shwartz and Tewari 2009
- Other common technique = LARS
 - □ Least angle regression and shrinkage, Efron et al. 2004

Soft Thresholding
$$\hat{w}_{\ell} = \left\{ \begin{array}{cc} (c_{\ell} + \lambda)/a_{\ell} & c_{\ell} < -\lambda \\ 0 & c_{\ell} \in [-\lambda, \lambda] \\ (c_{\ell} - \lambda)/a_{\ell} & c_{\ell} > \lambda \end{array} \right.$$
 From Kevin Murphy textbook









What you need to know



- Variable Selection: find a sparse solution to learning problem
- L₁ 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

©2005-2013 Carlos Guestrin