

#### Minimizing the Ridge Regression Objective



$$\hat{\mathbf{w}}_{ridge} = \arg\min_{w} \sum_{j=1}^{N} \left( t(x_j) - (w_0 + \sum_{i=1}^{k} w_i h_i(x_j)) \right)^2 + \lambda \sum_{i=1}^{k} w_i^2$$
$$= (H\mathbf{w} - \mathbf{t})^T (H\mathbf{w} - \mathbf{t}) + \lambda \mathbf{w}^T I_{0+k} \mathbf{w}$$

©2005-2013 Carlos Guestria

-

#### **Shrinkage Properties**

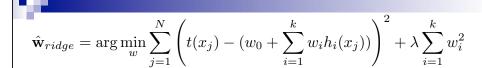


$$\hat{\mathbf{w}}_{ridge} = (H^T H + \lambda \ I_{0+k})^{-1} H^T \mathbf{t}$$

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

©2005-2013 Carlos Guestrin

#### Ridge Regression: Effect of Regularization

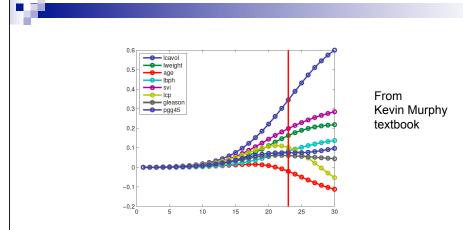


- Solution is indexed by the regularization parameter λ
- Larger λ
- Smaller λ
- As  $\lambda \rightarrow 0$
- As λ →∞

©2005-2013 Carlos Guestrin

7



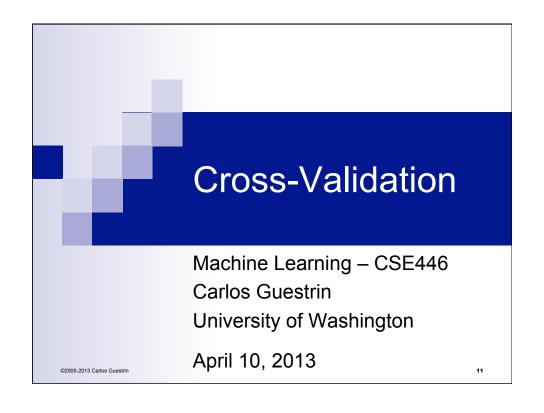


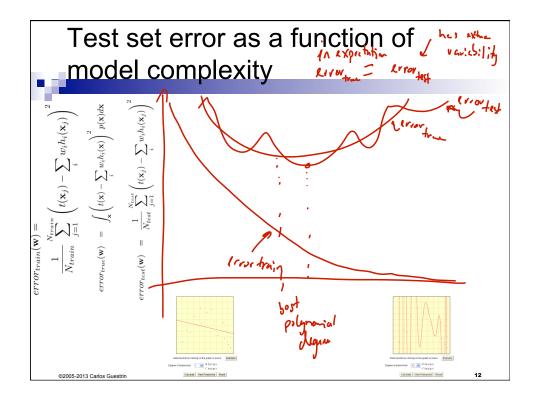
 Typical approach: select λ using cross validation, more on this later in the quarter

©2005-2013 Carlos Guestr

Error as a function of regularization parameter for a fixed model complexity 
$$\sum_{x_{luniv}} \frac{1}{\sum_{x_{luniv}} \frac{1}{\sum_{x_{lun$$

# What you need to know... Regularization Penalizes for complex models Ridge regression L<sub>2</sub> penalized least-squares regression Regularization parameter trades off model complexity with training error





#### How... 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... 🕾
- We could use the test data, but...

#### (LOO) Leave-one-out cross validation



- Consider a validation set with 1 example:
  - □ D training data
  - $\Box$  D\j training data with j th data point moved to validation set
- Learn classifier h<sub>D\i</sub> with D\i dataset
- **Estimate true error** as squared error on predicting t(x<sub>i</sub>):
  - □ Unbiased estimate of error<sub>true</sub>( $\boldsymbol{h}_{D_i}$ )!
  - □ Seems really bad estimator, but wait!
- LOO cross validation: Average over all data points *j*:
  - $\Box$  For each data point you leave out, learn a new classifier  $h_{Di}$
  - Estimate error as:

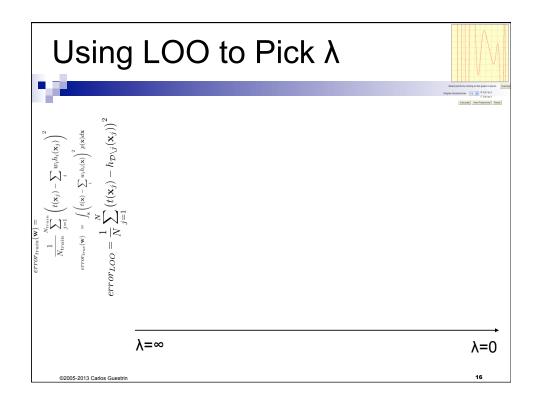
Estimate error as: 
$$error_{LOO} = rac{1}{N} \sum_{j=1}^{N} ig(t(\mathbf{x}_j) - h_{\mathcal{D}\setminus j}(\mathbf{x}_j)ig)^2$$

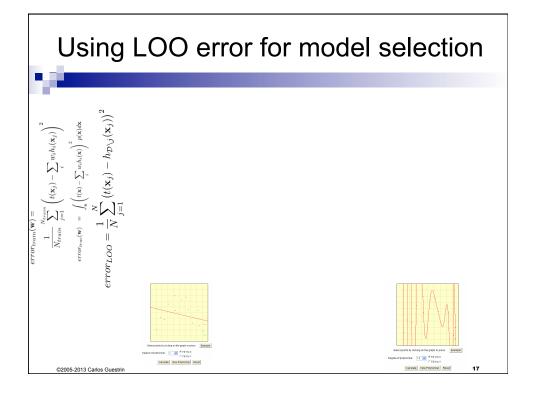
## LOO cross validation is (almost) unbiased estimate of true error of $h_D$ !

- 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 Guestria





#### 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 foooooreeeve'!!!
- Solution 1: Preferred, but not usually possible
  - ☐ Find a cool trick to compute LOO (e.g., see homework)

©2005-2013 Carlos Guestrin

### 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  $h_{D \mid D_i}$  using data point not in  $D_i$
  - □ Estimate error of  $h_{D \setminus Di}$  on validation set  $D_i$ :

$$error_{\mathcal{D}_i} = \frac{k}{N} \sum_{\mathbf{x}_j \in \mathcal{D}_i} (t(\mathbf{x}_j) - h_{\mathcal{D} \setminus \mathcal{D}_i}(\mathbf{x}_j))^2$$

• 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
  - $\square$  More (pessimistically) biased using much less data, only m(k-1)/k
  - ☐ Usually, **k = 10** ③

©2005-2013 Carlos Guestria

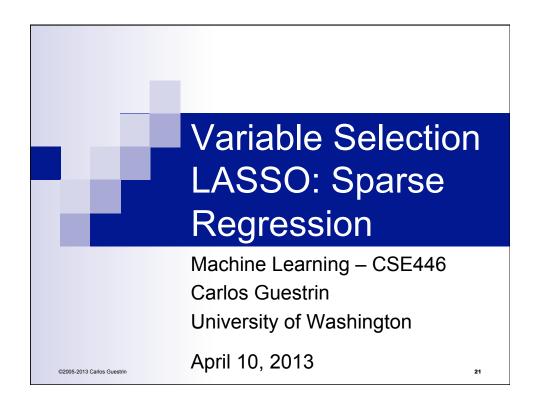
19

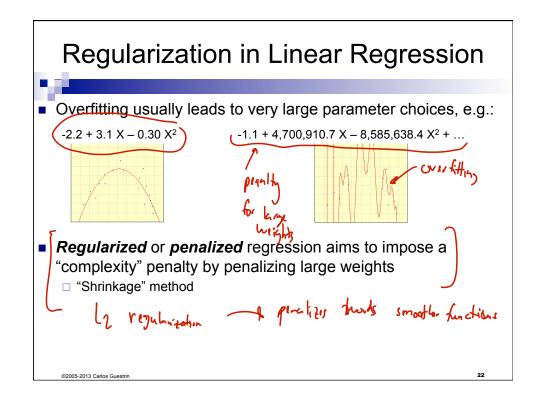
#### What you need to know...



- 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 Guestri





#### Variable Selection



- 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
  - □ Computationally intractable to perform "all subsets" regression
- Try new penalty: Penalize non-zero weights
  - □ Regularization penalty:
  - Leads to sparse solutions
  - $\ \square$  Just like ridge regression, solution is indexed by a continuous param  $\lambda$
  - □ This simple approach has changed statistics, machine learning & electrical engineering

©2005-2013 Carlos Guestrin

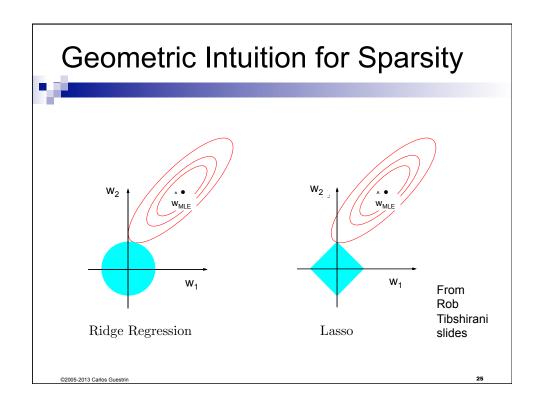
23

#### 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



- 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
  - □ Converges to optimum in some cases, such as LASSO

2

## Optimizing LASSO Objective One Coordinate at a Time



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

- 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:

©2005-2013 Carlos Guestrin

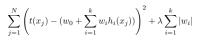
#### **Subgradients of Convex Functions**



Gradients lower bound convex functions:

- Gradients are unique at w iff function differentiable at w
- Subgradients: Generalize gradients to non-differentiable points:
  - ☐ Any plane that lower bounds function:

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





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

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

Gradient of RSS term: 
$$a_{\ell} = 2\sum_{j=1}^{N}(h_{\ell}(\mathbf{x}_{j}))^{2}$$
 
$$\frac{\partial}{\partial w_{\ell}}RSS(\mathbf{w}) = a_{\ell}w_{\ell} - c_{\ell}$$
 
$$c_{\ell} = 2\sum_{j=1}^{N}h_{\ell}(\mathbf{x}_{j})\left(t(\mathbf{x}_{j}) - (w_{0} + \sum_{i \neq \ell}w_{i}h_{i}(\mathbf{x}_{j}))\right)$$

- □ If no penalty:
- Subgradient of full objective:

#### Setting Subgradient to 0



$$\partial_{w_{\ell}} F(\mathbf{w}) = \begin{cases} a_{\ell} w_{\ell} - c_{\ell} - \lambda & w_{\ell} < 0 \\ [-c_{\ell} - \lambda, -c_{\ell} + \lambda] & w_{\ell} = 0 \\ a_{\ell} w_{\ell} - c_{\ell} + \lambda & w_{\ell} > 0 \end{cases}$$

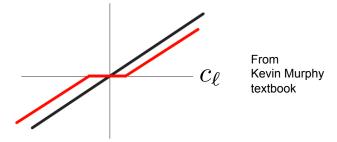
©2005-2013 Carlos Guestria

31

#### Soft Thresholding

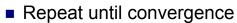


$$\hat{w}_{\ell} = \begin{cases} (c_{\ell} + \lambda)/a_{\ell} & c_{\ell} < -\lambda \\ 0 & c_{\ell} \in [-\lambda, \lambda] \\ (c_{\ell} - \lambda)/a_{\ell} & c_{\ell} > \lambda \end{cases}$$



©2005-2013 Carlos Guestrin

## Coordinate Descent for LASSO (aka Shooting Algorithm)



□ Pick a coordinate *l* at (random or sequentially)

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

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

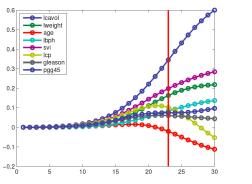
 $\hfill\Box$  For convergence rates, see Shalev-Shwartz and Tewari 2009

Other common technique = LARS

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

Suestrin

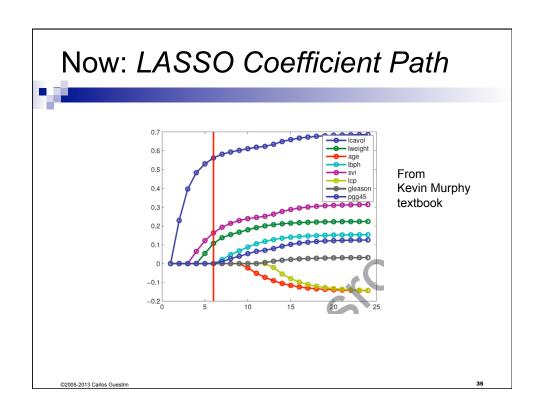
# Recall: Ridge Coefficient Path



From Kevin Murphy textbook

Typical approach: select λ using cross validation

©2005-2013 Carlos Guestri



LAS	SO Ex	ample			
	Term	Least Squares	Ridge	Lasso	
	Intercept	2.465	2.452	2.468	
	lcavol	0.680	0.420	0.533	From
	lweight	0.263	0.238	0.169	Rob Tibshirani
	age	-0.141	-0.046		slides
	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		
©2005-2013 Carlos Gi	uestrin				36

#### What you need to know



- Variable Selection: find a sparse solution to learning problem
- L<sub>1</sub> 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