

Regularization, Ridge Regression

Machine Learning – CSEP546

Carlos Guestrin

University of Washington

January 13, 2014

The regression problem

$$h_3(x) = \sin x \ln x$$

$$h_4(x) = e^{\cos x}$$

- Instances: $\langle \mathbf{x}_j, t_j \rangle$

- Learn: Mapping from x to $t(x)$

- Hypothesis space:

- Given, basis functions
- Find coeffs $\mathbf{w} = \{w_1, \dots, w_k\}$

- Why is this called linear regression???

- model is linear in the parameters

$$H = \{h_1, \dots, h_K\}$$

$$t(\mathbf{x}) \approx \hat{f}(\mathbf{x}) = \sum_i w_i h_i(\mathbf{x})$$

linear combination

- Precisely, minimize the residual squared error:

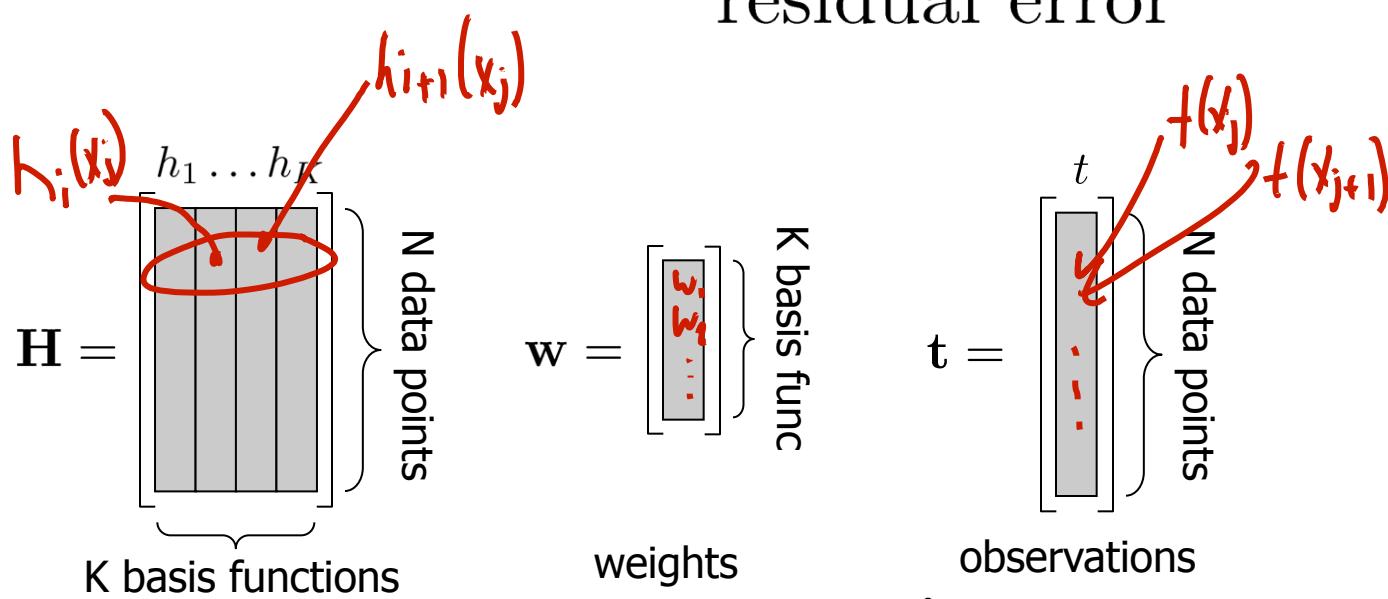
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_{j=1}^N \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

N *sum data*
 truth
 estimate
 squared

The regression problem in matrix notation

$$\cancel{\mathbf{w}^*} = \arg \min_{\mathbf{w}} \sum_{j=1}^N \left(t(\mathbf{x}_j) - \sum_{i=1}^K w_i h_i(\mathbf{x}_j) \right)^2$$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \underbrace{(\mathbf{H}\mathbf{w} - \mathbf{t})^T (\mathbf{H}\mathbf{w} - \mathbf{t})}_{\text{residual error}}$$



Regression solution = simple matrix operations

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \underbrace{(\mathbf{H}\mathbf{w} - \mathbf{t})^T (\mathbf{H}\mathbf{w} - \mathbf{t})}_{\text{residual error}} \quad \text{loss function}$$

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

where $\mathbf{A} = \mathbf{H}^T \mathbf{H} = \begin{bmatrix} & & \\ & & \\ & & \\ & & \\ & & \end{bmatrix}$ $\mathbf{b} = \mathbf{H}^T \mathbf{t} = \begin{bmatrix} & \\ & \\ & \\ & \end{bmatrix}$

$\underbrace{\qquad\qquad\qquad}_{k \times k \text{ matrix}}$ $\underbrace{\qquad\qquad\qquad}_{k \times 1 \text{ vector}}$

for k basis functions

Bias-Variance Tradeoff

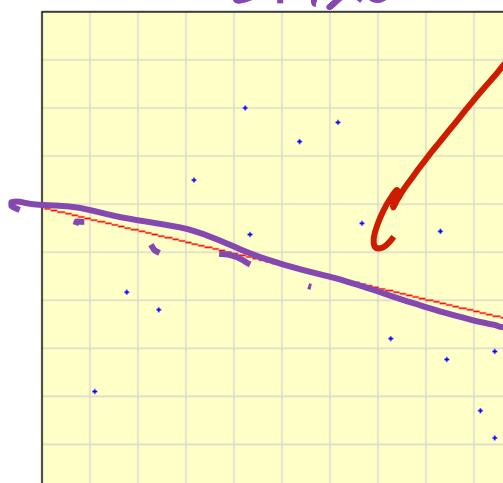
e.g. degree of polynomial

- Choice of hypothesis class introduces learning bias

□ More complex class → less bias

□ More complex class → more variance *h_D very diff from h_{D'}*

Lagged



Select points by clicking on the graph or press

Example

Degree of polynomial: Fit Y to X
 Fit X to Y

Calculate View Polynomial Reset

in between

Select points by clicking on the graph or press

Example

Degree of polynomial: Fit Y to X
 Fit X to Y

Calculate View Polynomial Reset

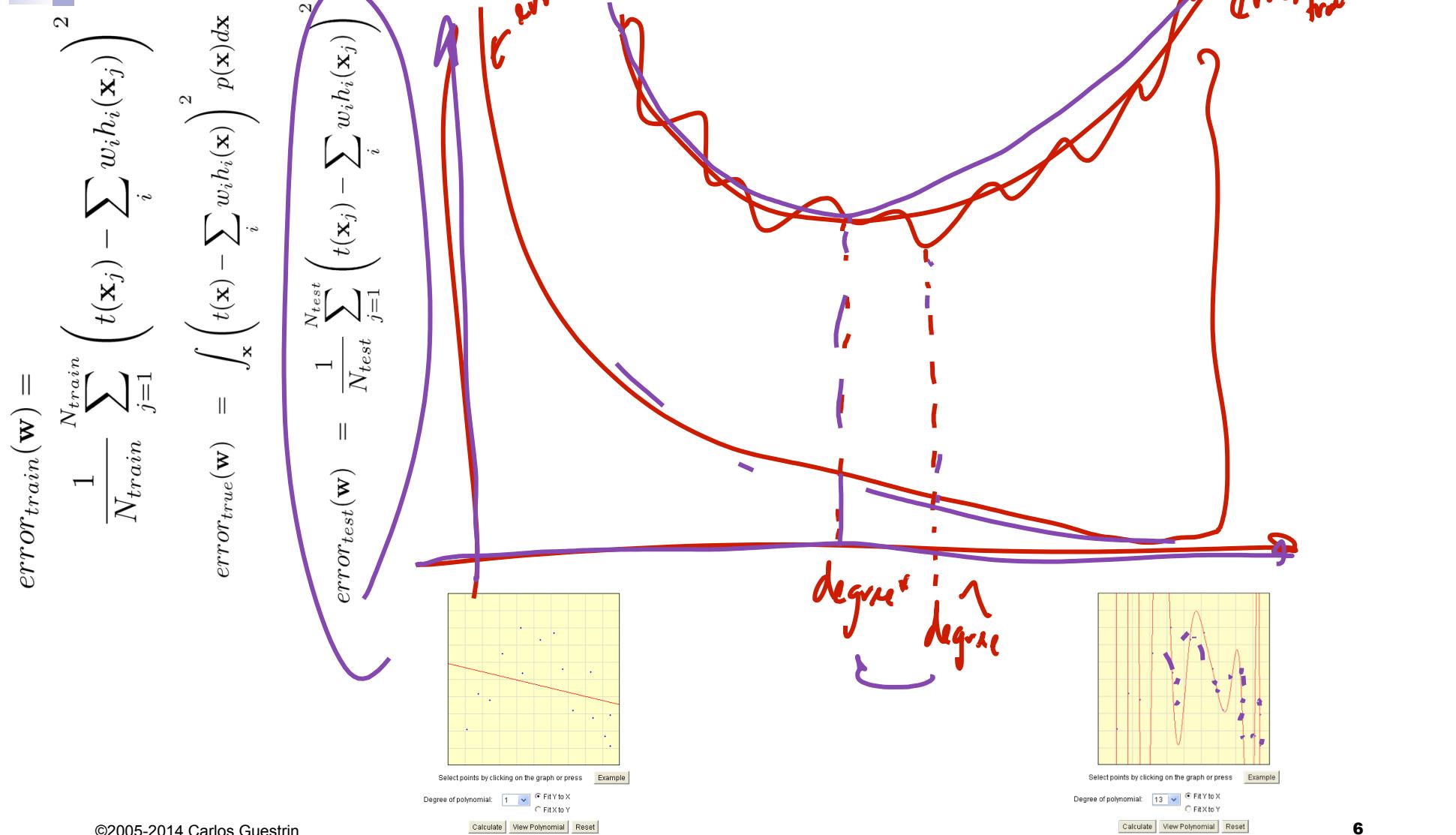
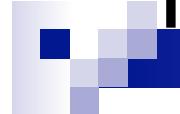
Select points by clicking on the graph or press

Example

Degree of polynomial: Fit Y to X
 Fit X to Y

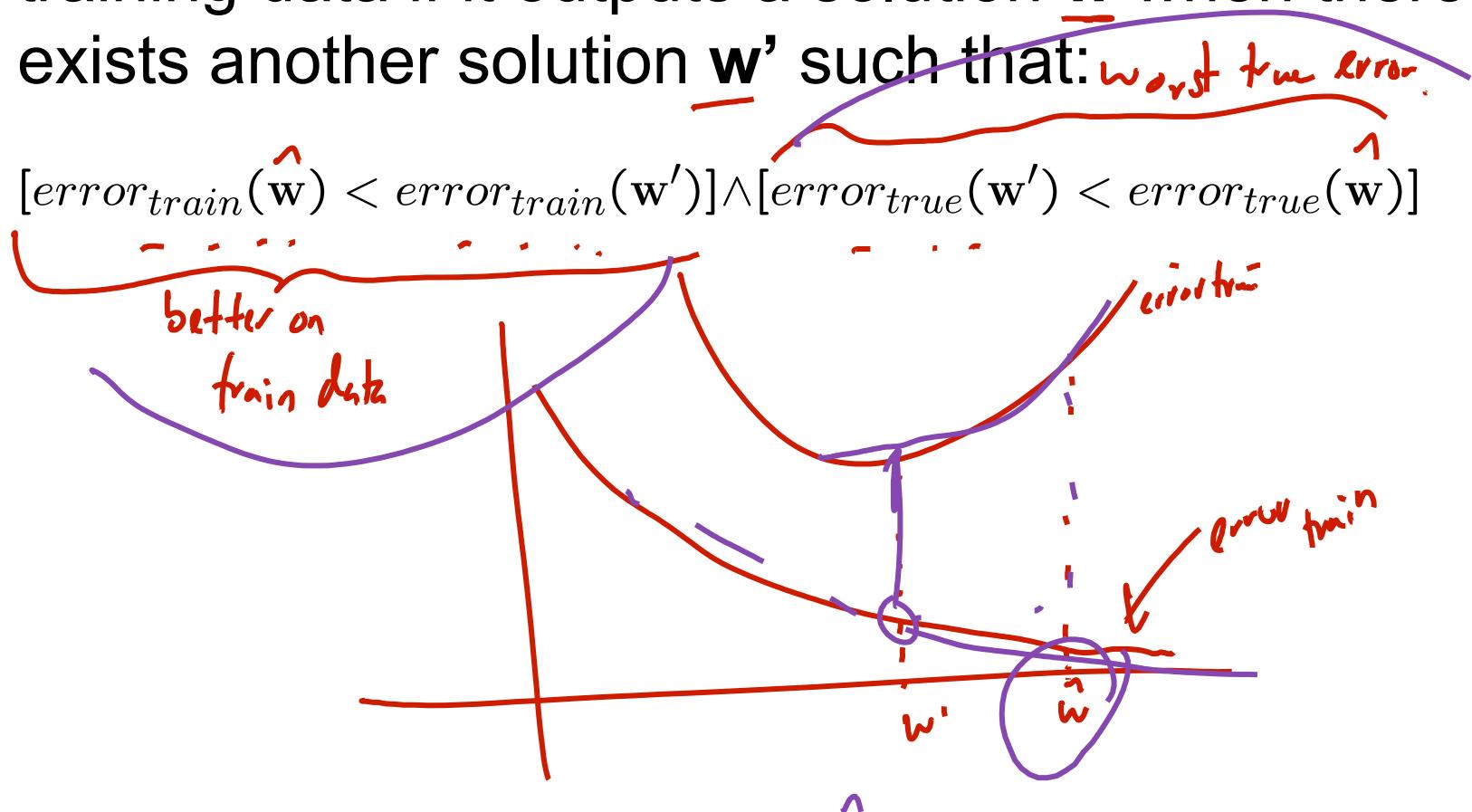
Calculate View Polynomial Reset

Test set error as a function of model complexity



Overfitting

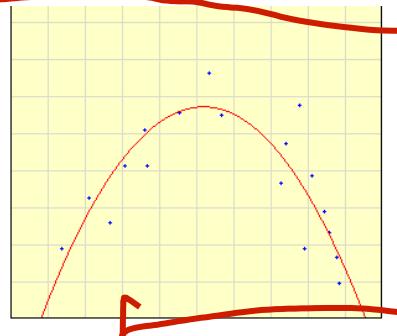
- **Overfitting:** a learning algorithm overfits the training data if it outputs a solution \hat{w} when there exists another solution w' such that: *worst true error*.



Regularization in Linear Regression

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

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



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



sweet spot

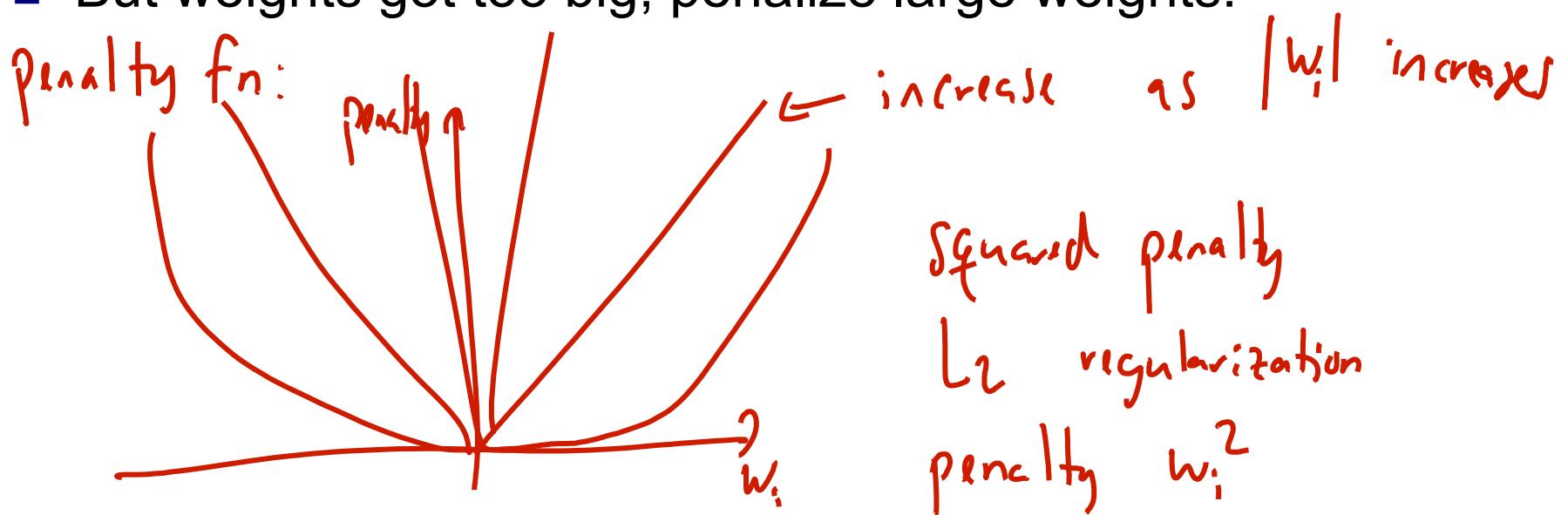
- Regularized or penalized** regression aims to impose a “complexity” penalty by penalizing large weights
 - “Shrinkage” method

Quadratic Penalty (regularization)

- What we thought we wanted to minimize:

Squared error in prediction

- But weights got too big, penalize large weights:



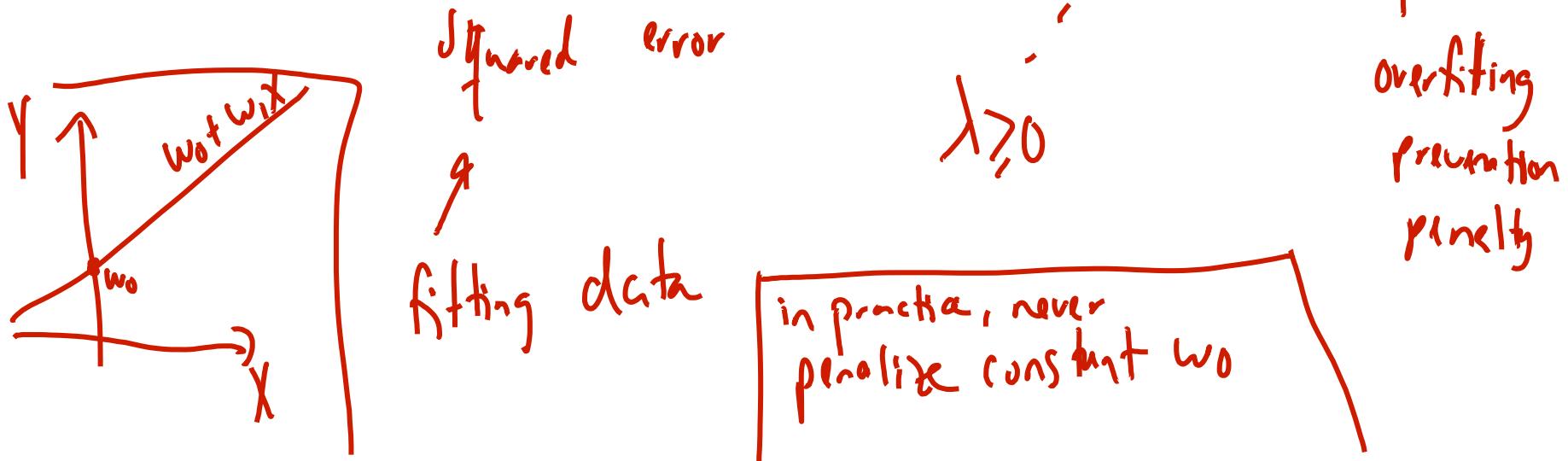
Ridge Regression

$$\|w\|_2^2 = \sum_{i=1}^k w_i^2$$

$w = (w_1 \dots w_k)$

- Ameliorating issues with overfitting: *penalizing large weights*
- New objective:

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



Ridge Regression in Matrix Notation

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

$$\mathbf{H} = \begin{bmatrix} h_0 & \dots & h_k \end{bmatrix} \quad \begin{array}{c} \text{K+1 basis functions} \\ \underbrace{\hspace{10em}}_{\text{N data points}} \end{array}$$
$$\mathbf{w} = \begin{bmatrix} \text{weights} \end{bmatrix} \quad \begin{array}{c} \text{K+1 basis func} \\ \underbrace{\hspace{10em}}_{\text{N data points}} \end{array}$$
$$\mathbf{t} = \begin{bmatrix} \text{observations} \end{bmatrix} \quad \begin{array}{c} \text{Same as regression} \\ \text{N data points} \end{array}$$

$$I_{0+k} = \begin{pmatrix} 0 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & 1 & \\ & & & & 0 \\ & & & & & 1 \\ & & & & & & 1 \\ & & & & & & & 1 \end{pmatrix}$$

Minimizing the Ridge Regression Objective

$$\frac{1}{n} \text{ as and } \hat{\mathbf{w}}_{\text{MLE}}^* = \underbrace{(\mathbf{H}^T \mathbf{H})^{-1}}_{\mathbf{A}^{-1}} \underbrace{\mathbf{H}^T \mathbf{t}}_{\mathbf{b}} = \mathbf{A}^{-1} \mathbf{b}$$

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

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

Shrinkage Properties

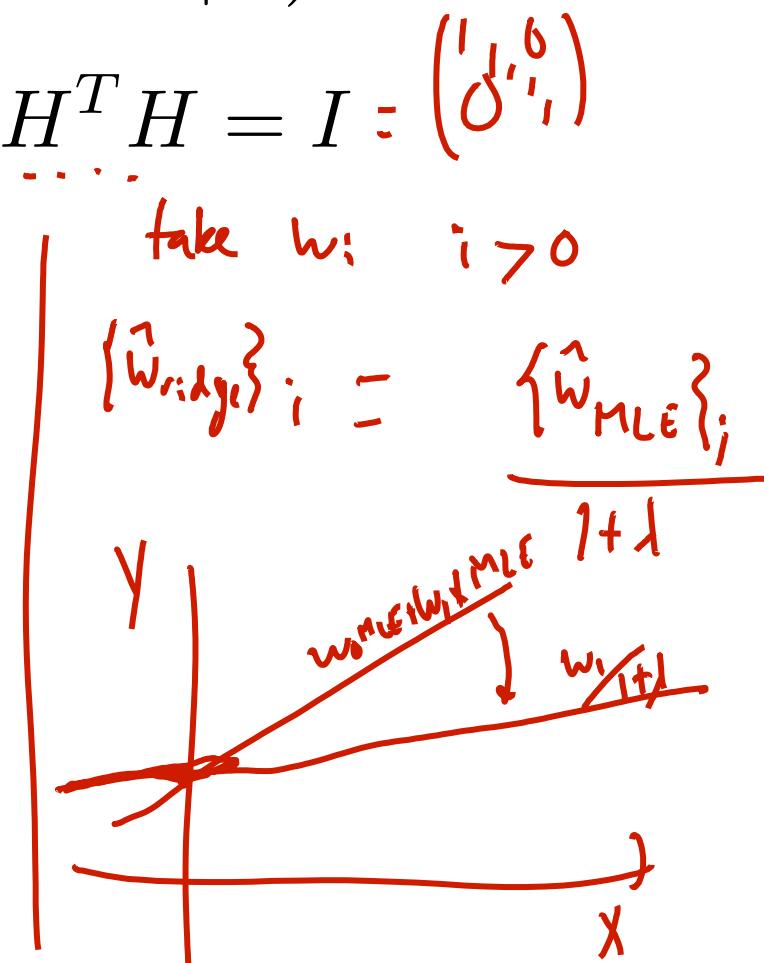
$$\hat{\mathbf{w}}_{ridge} = (H^T H + \lambda I_{0+k})^{-1} H^T \mathbf{t}$$

- If orthonormal features/basis: $H^T H = I$: $\begin{pmatrix} 1 & 0 \\ 0 & \ddots \\ 0 & 0 \end{pmatrix}$

$$\hat{\mathbf{w}}_{ridge} = \underbrace{(I + \lambda I_{0+k})^{-1}}_{\begin{pmatrix} 1 & 0 \\ 0 & \ddots \\ 0 & 0 \end{pmatrix}} H^T \mathbf{t}$$

$$\begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & \frac{1}{1+\lambda} & \dots & 0 \\ 0 & 0 & \ddots & 0 \\ \vdots & & & \frac{1}{1+\lambda} \end{pmatrix}$$

$$\hat{\mathbf{w}}_{MLE} = I^{-1} H^T \mathbf{t}$$



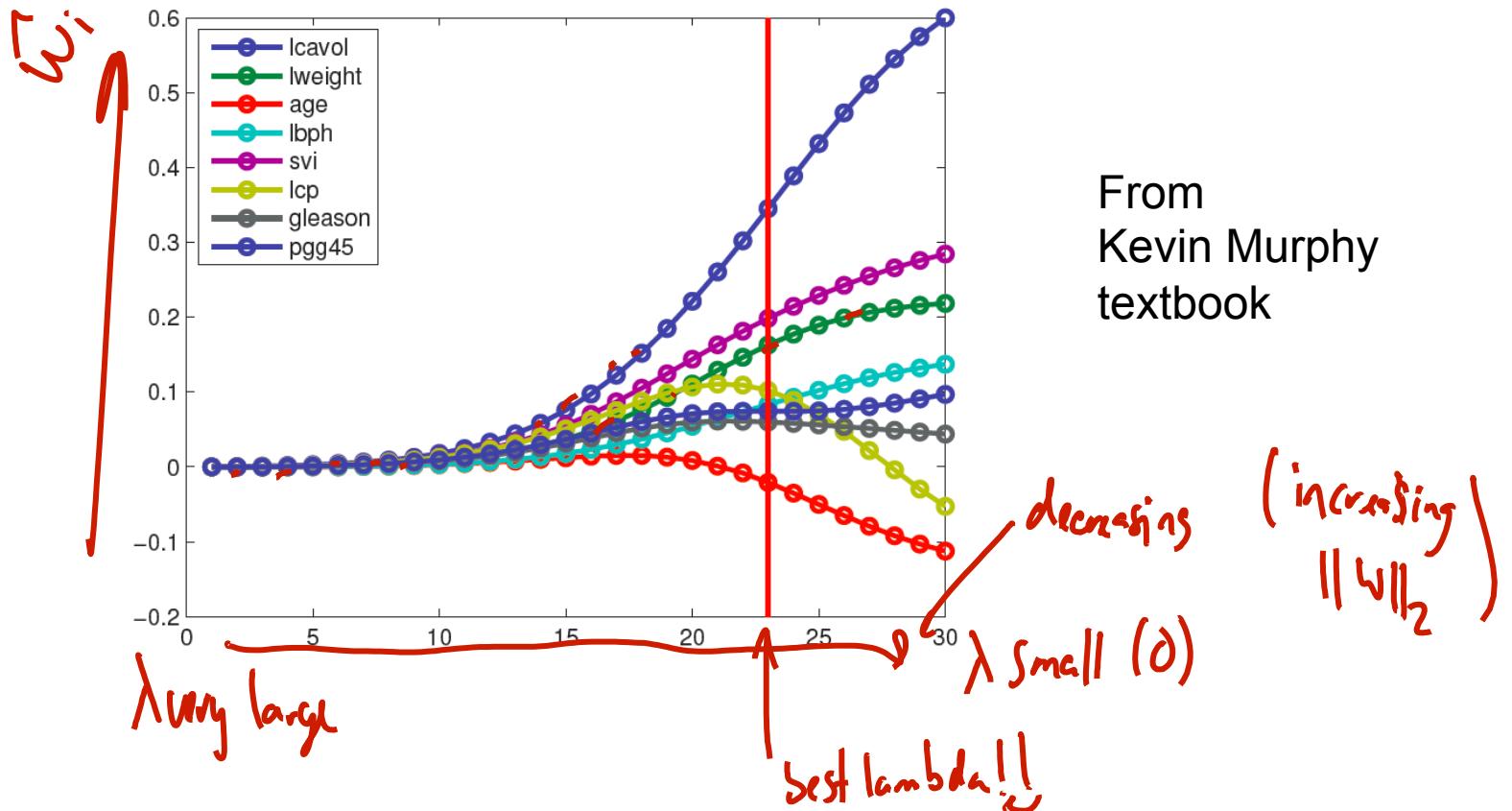
Ridge Regression: Effect of Regularization

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

- Solution is indexed by the regularization parameter λ
- Larger λ more penalty \rightarrow smoother , more bias
- Smaller λ less penalty \rightarrow more flexible , more variance
- As $\lambda \rightarrow 0$ $\hat{\mathbf{w}}_{ridge} \rightarrow \hat{\mathbf{w}}_{MLE}$
- As $\lambda \rightarrow \infty$ $w_1 \dots w_k \rightarrow 0$
get $w_0 \leftarrow \text{average } t(x_j)$

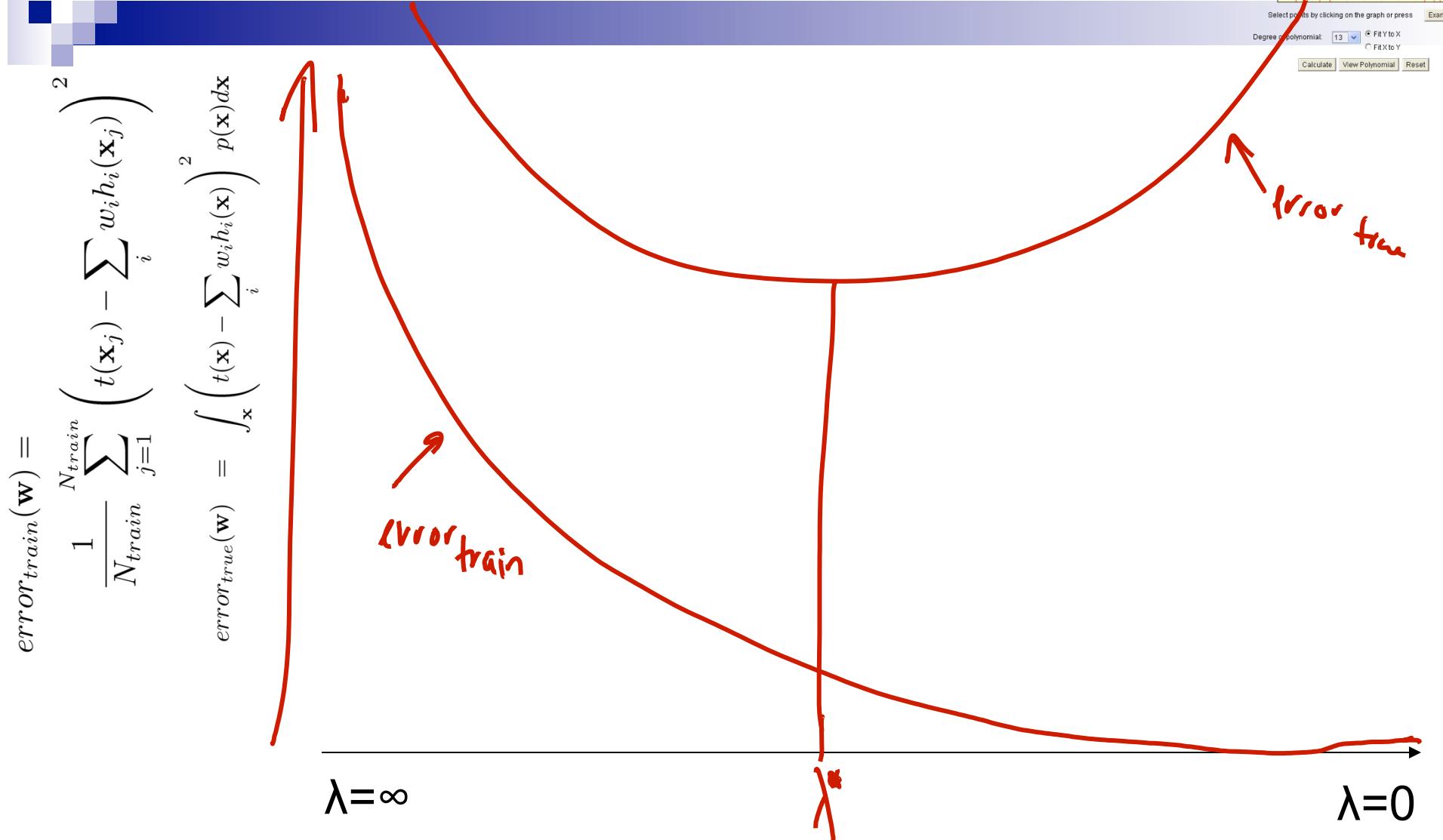
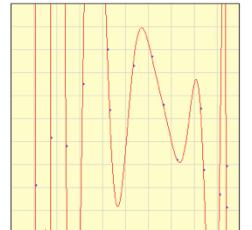
Ridge Coefficient Path

$$\|v\|_2 = \sqrt{\sum_i v_i^2}$$



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

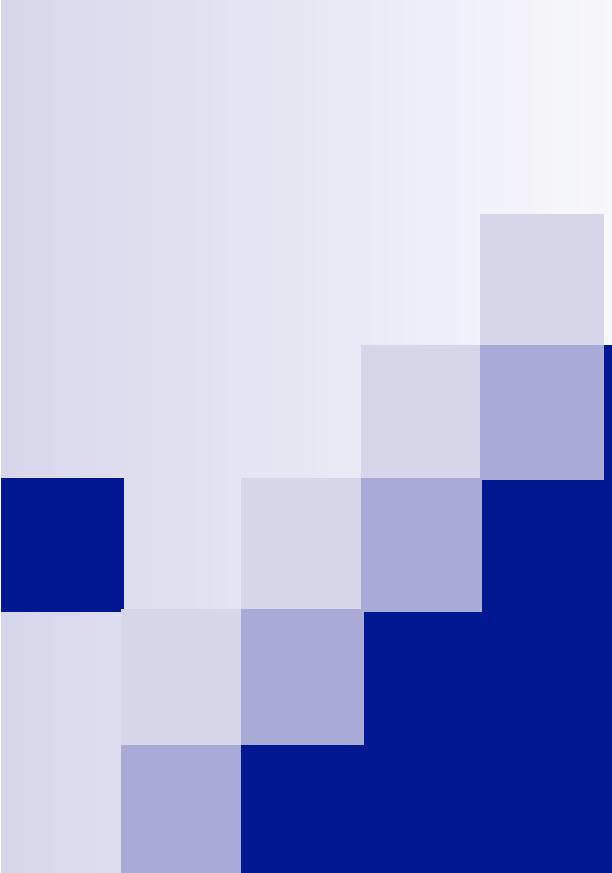
Error as a function of regularization parameter for a fixed model complexity



What you need to know...



- Regularization
 - Penalizes for complex models
- Ridge regression
 - L_2 penalized least-squares regression
 - Regularization parameter trades off model complexity with training error

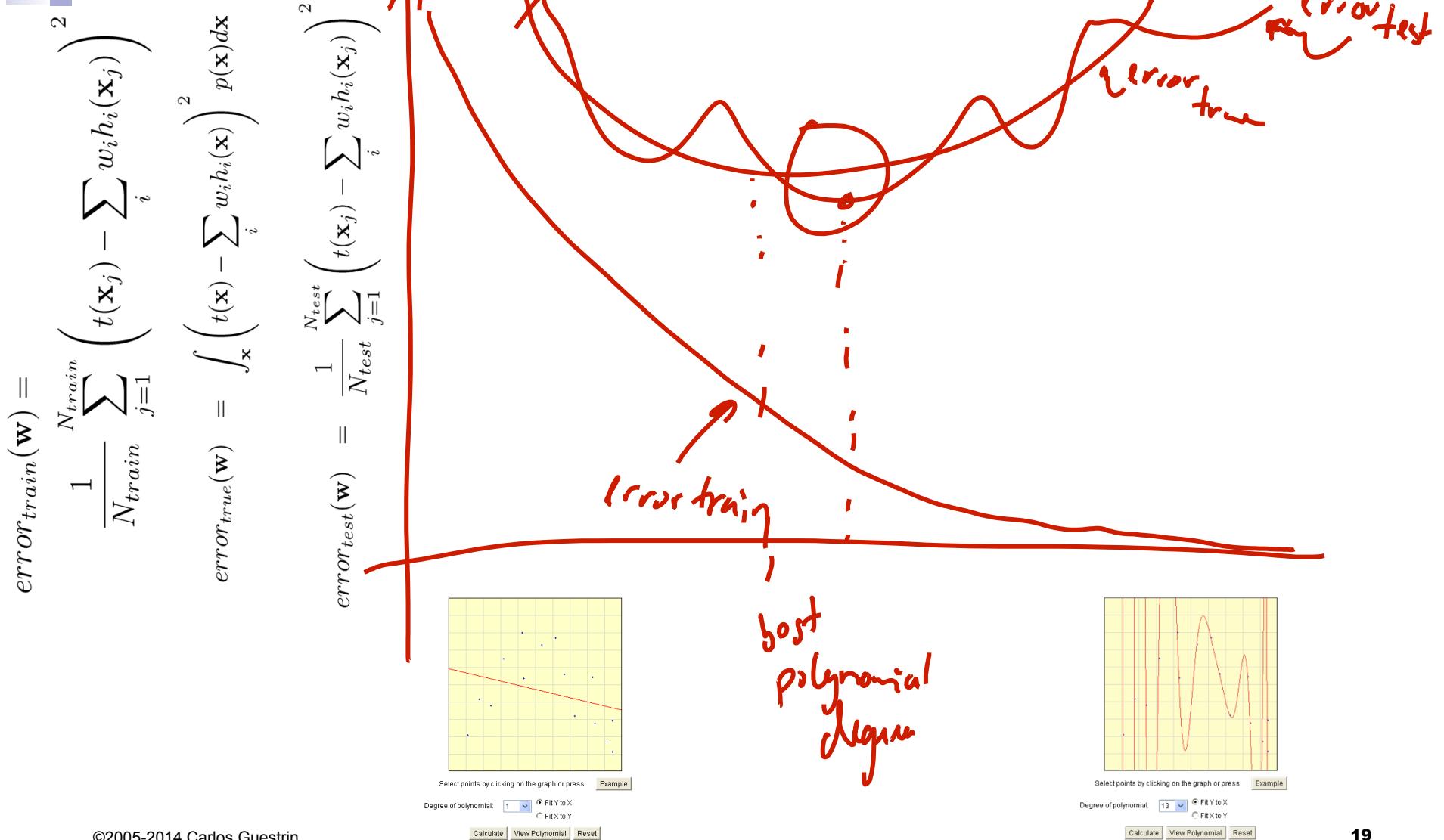
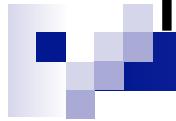


Cross-Validation

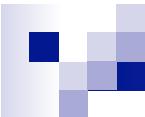
Machine Learning – CSEP546
Carlos Guestrin
University of Washington

January 13, 2014

Test set error as a function of model complexity



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...
 1. don't learn on test data
 2. don't have a lot of data

(LOO) Leave-one-out cross validation

- Consider a **validation set with 1 example:**

choose lambda based on
1 test example

- D – training data

- $D \setminus j$ – training data with j th data point moved to validation set

- Learn classifier $h_{D \setminus j}$ with $D \setminus j$ dataset \leftarrow do this for all j
regressor / learn N functions

- Estimate true error as squared error on predicting $t(\mathbf{x}_j)$:

- Unbiased estimate of error_{true}($h_{D \setminus j}$)!

$$\text{error}_{\text{true}}(h_{D \setminus j}) = E_{\mathbf{x}}[(t - h_{D \setminus j})^2]$$

- Seems really bad estimator, but wait!

- LOO cross validation: Average over all data points j :

- For each data point you leave out, learn a new classifier $h_{D \setminus j}$

- Estimate error as:

$$\text{error}_{\text{LOO}} = \frac{1}{N} \sum_{j=1}^N (t(\mathbf{x}_j) - h_{D \setminus j}(\mathbf{x}_j))^2$$

$\lambda \leftarrow$ parameterized by y

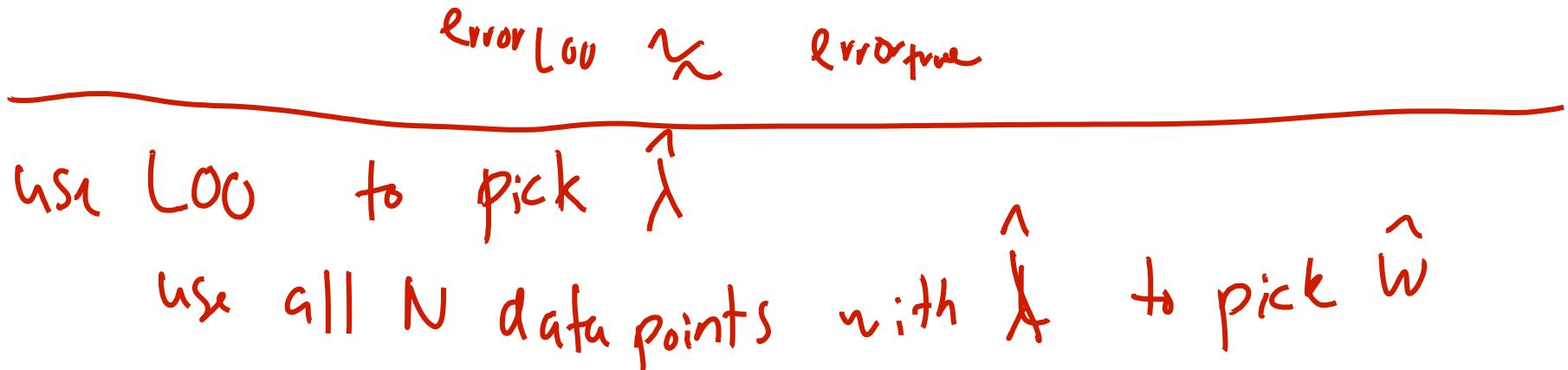
LOO cross validation is (almost) unbiased estimate of true error of h_D !

- 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!

$\text{error}_{\text{LOO}} \approx \text{error}_{\text{true}}$

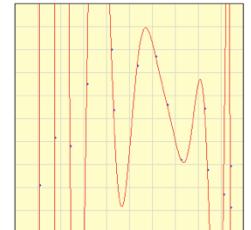
use LOO to pick $\hat{\lambda}$

use all N data points with $\hat{\lambda}$ to pick \hat{w}



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

Using LOO to Pick λ

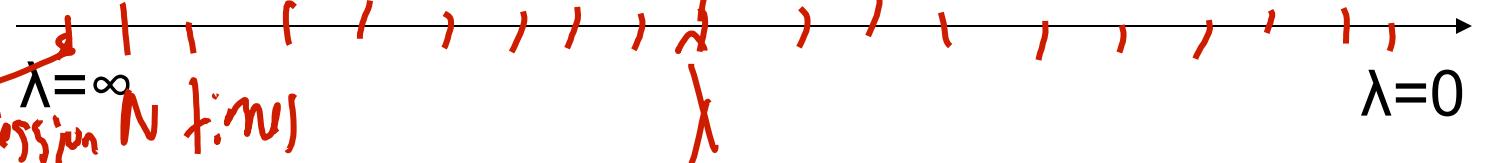


$$\text{error}_{\text{train}}(\mathbf{w}) = \frac{1}{N_{\text{train}}} \sum_{j=1}^{N_{\text{train}}} \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

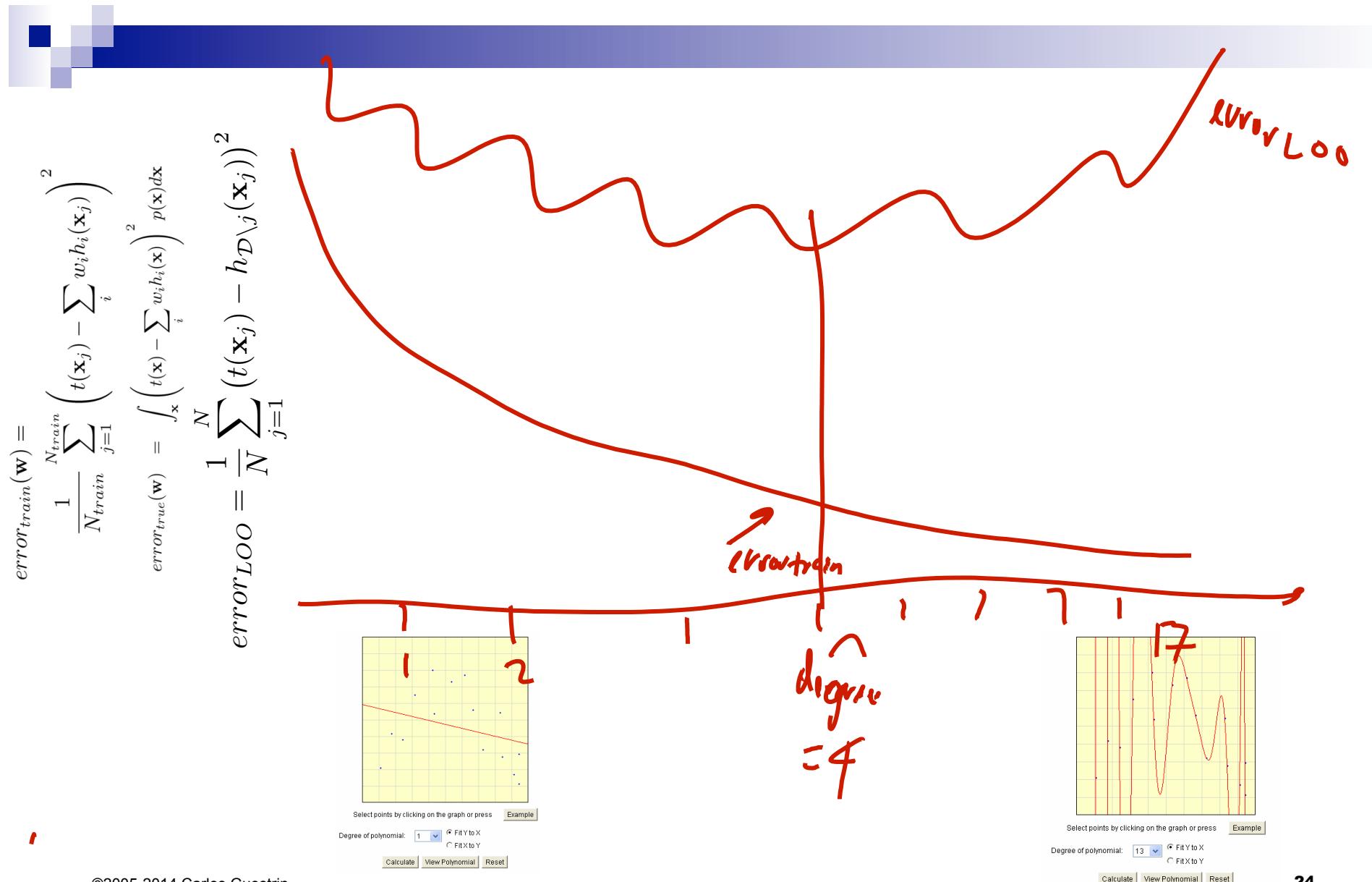
$$\text{error}_{\text{true}}(\mathbf{w}) = \int_{\mathbf{x}} \left(t(\mathbf{x}) - \sum_i w_i h_i(\mathbf{x}) \right)^2 p(\mathbf{x}) d\mathbf{x}$$

$$\text{error}_{\text{LOO}} = \frac{1}{N} \sum_{j=1}^N (t(\mathbf{x}_j) - h_{\mathcal{D} \setminus j}(\mathbf{x}_j))^2$$

$\lambda = \infty$
run regression N times



Using LOO error for model selection



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

extra credit

Solution 2 to complexity of computing LOO: (More typical) **Use k -fold cross validation**

- Randomly divide training data into k equal parts
 - D_1, \dots, D_k
- For each i
 - Learn classifier $h_{D \setminus D_i}$ using data point not in D_i
 - Estimate error of $h_{D \setminus D_i}$ on validation set D_i :
- **k -fold cross validation error is average** over data splits:

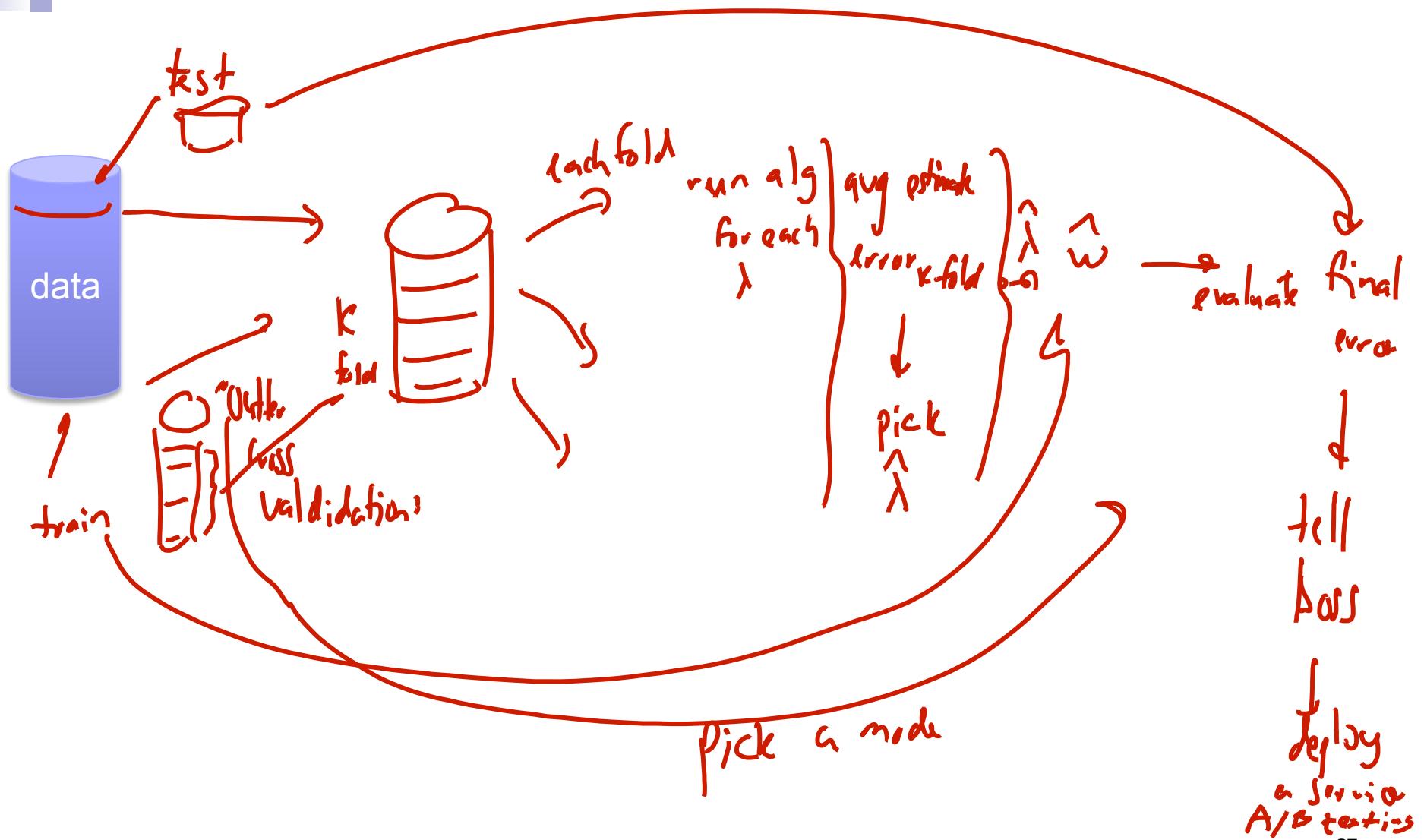
$$error_{\mathcal{D}_i} = \frac{k}{N} \sum_{x_j \in \mathcal{D}_i} (t(x_j) - h_{\mathcal{D} \setminus \mathcal{D}_i}(x_j))^2$$

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

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

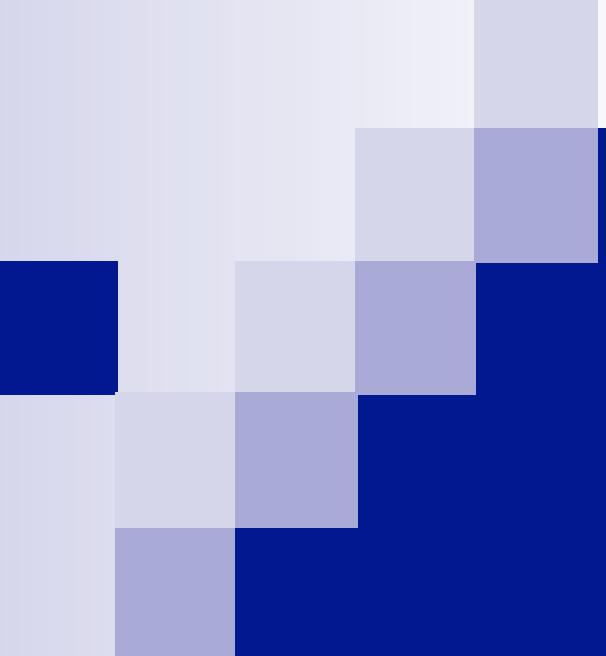


ML Pipeline



What you need to know...

- Never ever ever ever ever ever ever ever
ever ever ever ever ever ever ever ever
ever ever ever ever ever ever ever ever
train on the test data
- 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



Variable Selection LASSO: Sparse Regression

Machine Learning – CSEP546

Carlos Guestrin

University of Washington

January 13, 2014

Sparsity

age gender MC grade ..

GPA

- Vector w is sparse, if many entries are zero:

$$w = (0.8, 0, 0.0, 0.7, 0.0, -2, \dots)$$

- Very useful for many tasks, e.g.,
 - Efficiency:** If $\text{size}(w) = \underline{100B}$, each prediction is expensive:
 - If part of an online system, too slow
 - If w is sparse, prediction computation only depends on number of non-zeros
 - Interpretability:** What are the relevant dimension to make a prediction?
 - E.g., what are the parts of the brain associated with particular words?
- But computationally intractable to perform “all subsets” regression

(100B)
K
Subjects

Mean of independently learned signatures over all nine participants

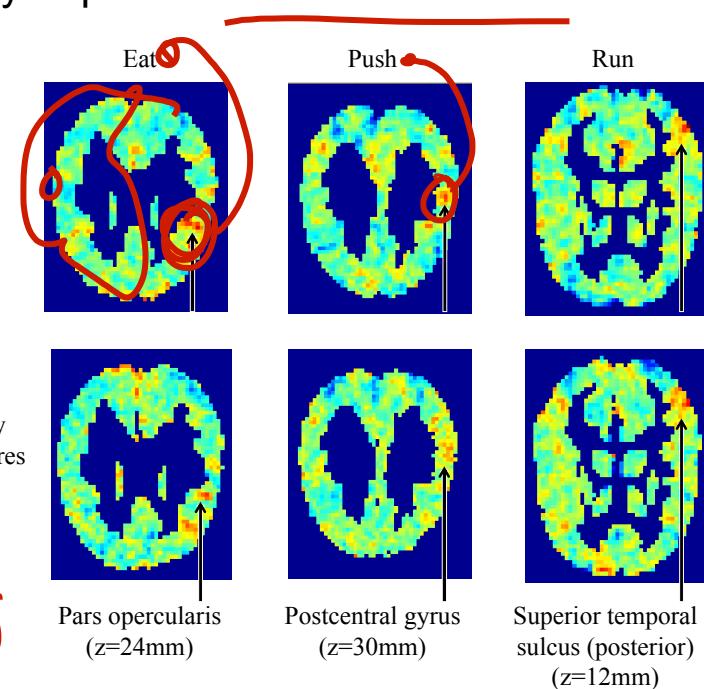


Figure from Tom Mitchell

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 $F_0 = \emptyset$
 - Run learning algorithm for current set of features F_t
 - Obtain h_t
 - Select **next best feature X_i^***
 - e.g., X_j that results in lowest training error learner when learning with $F_t + \{X_j\}$

decreases error the most
 - $F_{t+1} \leftarrow F_t + \{X_i^*\}$
 - Recurse

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

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

When do you stop???

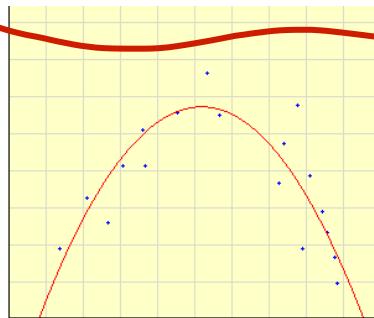
- When training error is low enough?
 - When test set error is low enough?
 -

Cross validation

Regularization in Linear Regression

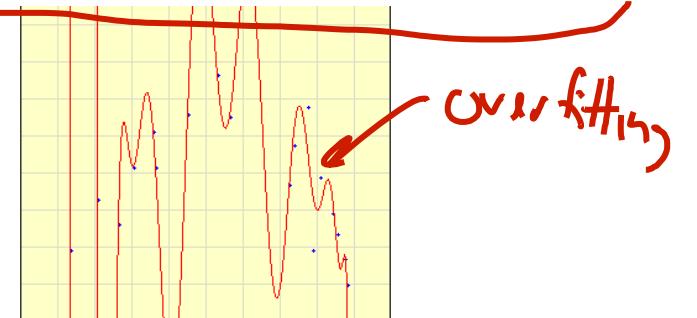
- Overfitting usually leads to very large parameter choices, e.g.:

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



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

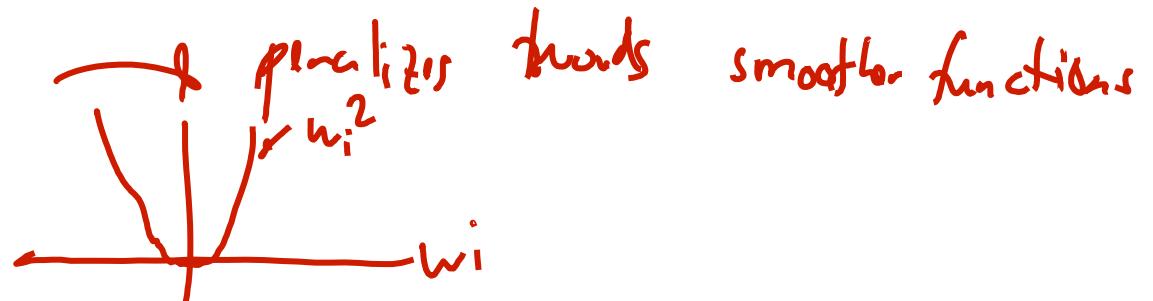
penalty
for large
weights



overfitting

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

L_2 regularization



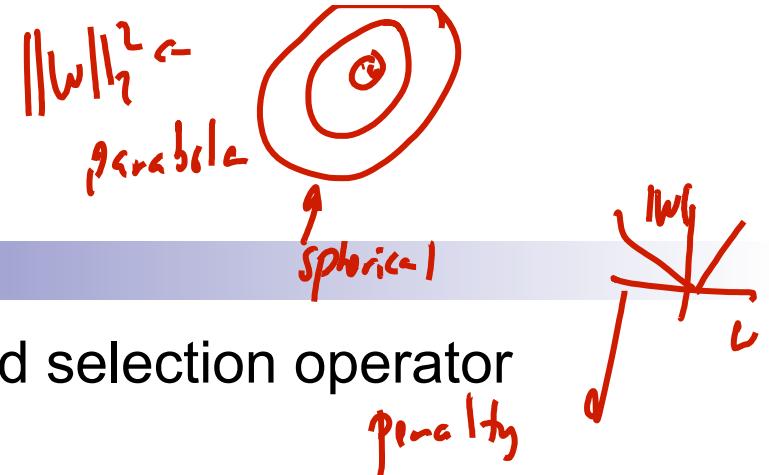
Variable Selection by Regularization

- 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

G-thresholding \rightarrow can lead to poor solution in practice
in ridge ($0.8 \dots -\epsilon + \epsilon, \epsilon - \epsilon -$)

- Try new penalty: Penalize non-zero weights
 - Regularization penalty:
instead $\lambda \|w\|_2^2 = \sum_i w_i^2$ → $\lambda \|w\|_1$
 - Leads to sparse solutions
 - Just like ridge regression, solution is indexed by a continuous param λ
 - This simple approach has changed statistics, machine learning & electrical engineering

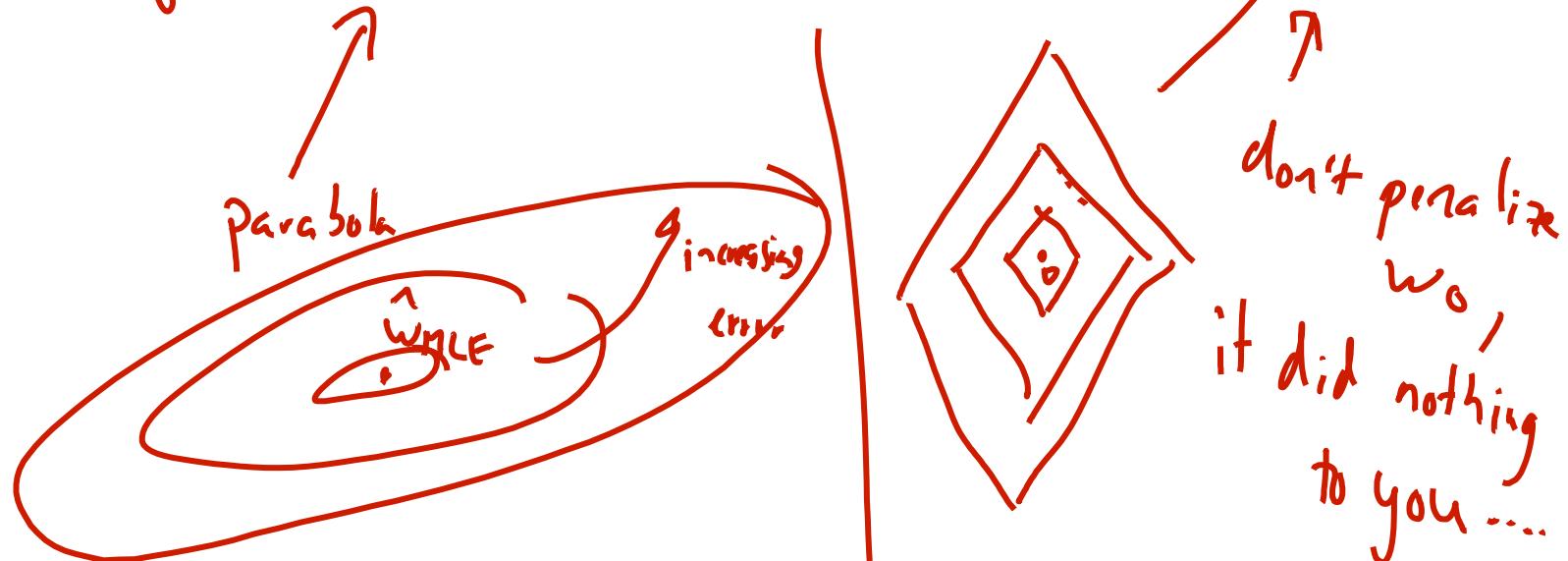
LASSO Regression



- LASSO: least absolute shrinkage and selection operator

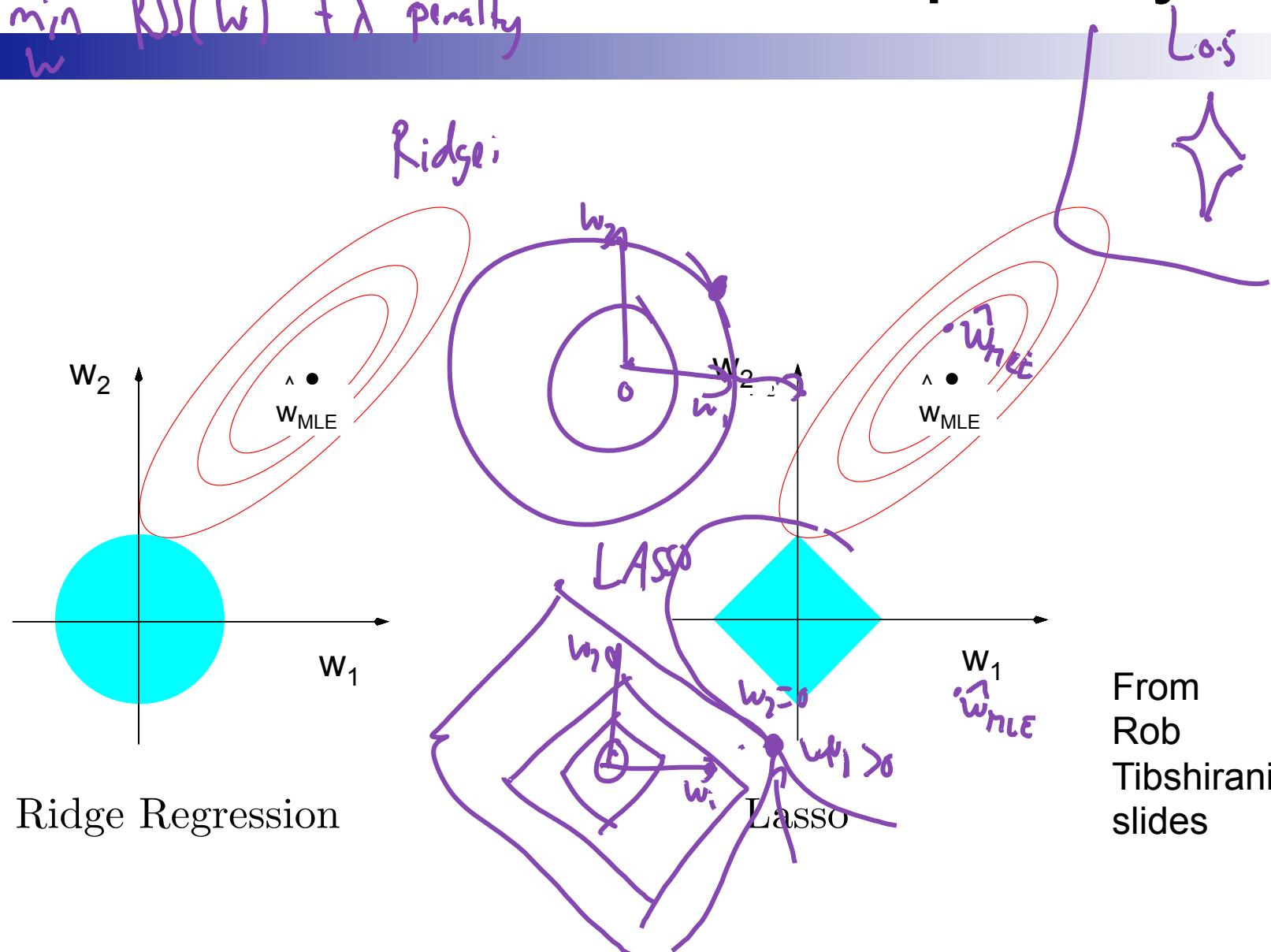
- New objective:

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



Geometric Intuition for Sparsity

$$\min_w \text{RSS}(w) + \lambda \text{ penalty}$$

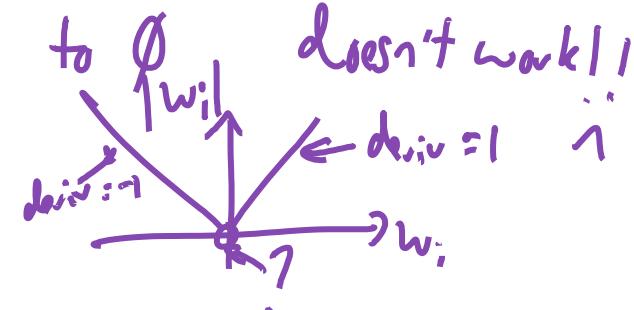


Optimizing the LASSO Objective

- LASSO solution:

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

Take derivative w.r.t w_i set to 0 doesn't work!!
1. derivative $|w_i|$



2. Even if you could take derivative (subgradient) there is no closed-form solution!!
for $\hat{\mathbf{w}}_{LASSO}$ in

Coordinate Descent

- Given a function F
 - Want to find minimum
$$\hat{w} = \underset{w}{\operatorname{argmin}} F(w_0, w_1, \dots, w_k)$$
- Often, hard to find minimum for all coordinates, but easy for one coordinate
1-d optimization problems
- Coordinate descent:
initialize $w = 0$ or something else
while not converged,
pick coordinate i
$$\hat{w}_i \leftarrow \underset{z}{\operatorname{argmin}} F(w_0, \hat{w}_1, \dots, \hat{w}_{i-1}, z, \hat{w}_{i+1}, \dots, \hat{w}_k)$$

fix all but w_i
optimize w_i
- How do we pick next coordinate?
round robin, random, "smartly"
- Super useful approach for *many* problems
 - Converges to optimum in some cases, such as LASSO

testing convergence:
round robin: diff in
error between passes
over all parameters
 $F(w^{(0)})$
diff pass over all params
same $F(w^{(1)})$ → Stop
small $\Rightarrow \epsilon$

LASSO OPT !!
Other times: local optima
or worse

How do we find the minimum over each coordinate?

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

$$\underset{z}{\operatorname{argmin}} \quad F(w_0, w_1, \dots w_{i-1}, z, w_{i+1}, \dots w_k)$$

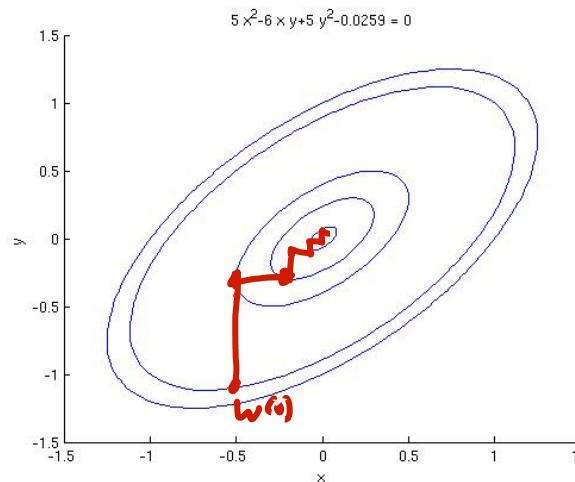


Illustration from Wikipedia

- Standard approach: *partial*
take derivative w.r.t w_i & set to 0

Optimizing LASSO Objective

One Coordinate at a Time

RSS(\mathbf{w})

penalty

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

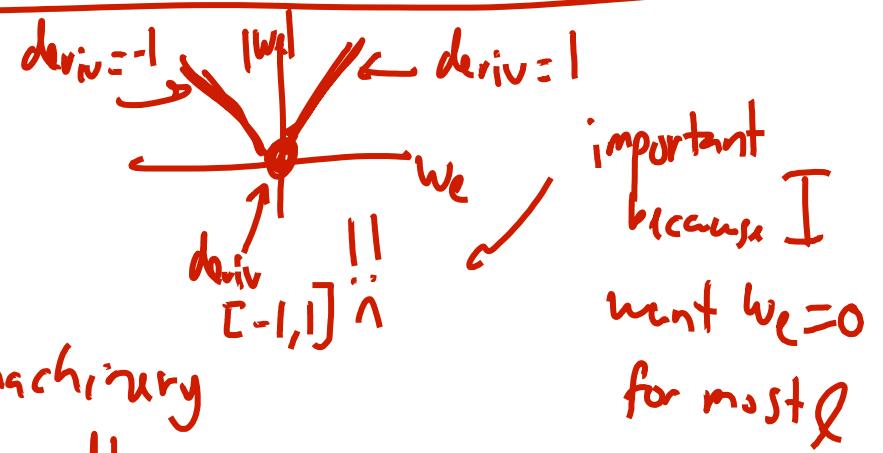
- Taking the derivative:

- 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:

$$\frac{d}{\partial w_\ell} \lambda \sum_{i=1}^k |w_i| = \lambda \frac{d}{\partial w_\ell} |w_\ell|$$



Need more advanced machinery
called subgradient

Coordinate Descent for LASSO (aka Shooting Algorithm)

initialize $w = 0 \leftarrow$ for first λ

Outer loop explores decreasing values of λ
using soln from one λ to init for next
stop $\text{old_error} - \text{new_error} < \epsilon$

old error = RSS(w) / Mult
new error ...

for $\ell = 0 \dots k$

- Repeat until convergence

- Pick a coordinate ℓ at (random or sequentially)

Hl 70

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

- Where:

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

recompute ↵

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

depends on all \hat{w}_i except for \hat{w}_ℓ

don't regularize w_0

$$w_0 = c_0/a_0$$

$$\hat{w}_0 = \frac{1}{N} \sum_{j=1}^N \left(t(\mathbf{x}_j) - \sum_{i=1}^k \hat{w}_i h_i(\mathbf{x}_j) \right)$$

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

- Other common technique = LARS

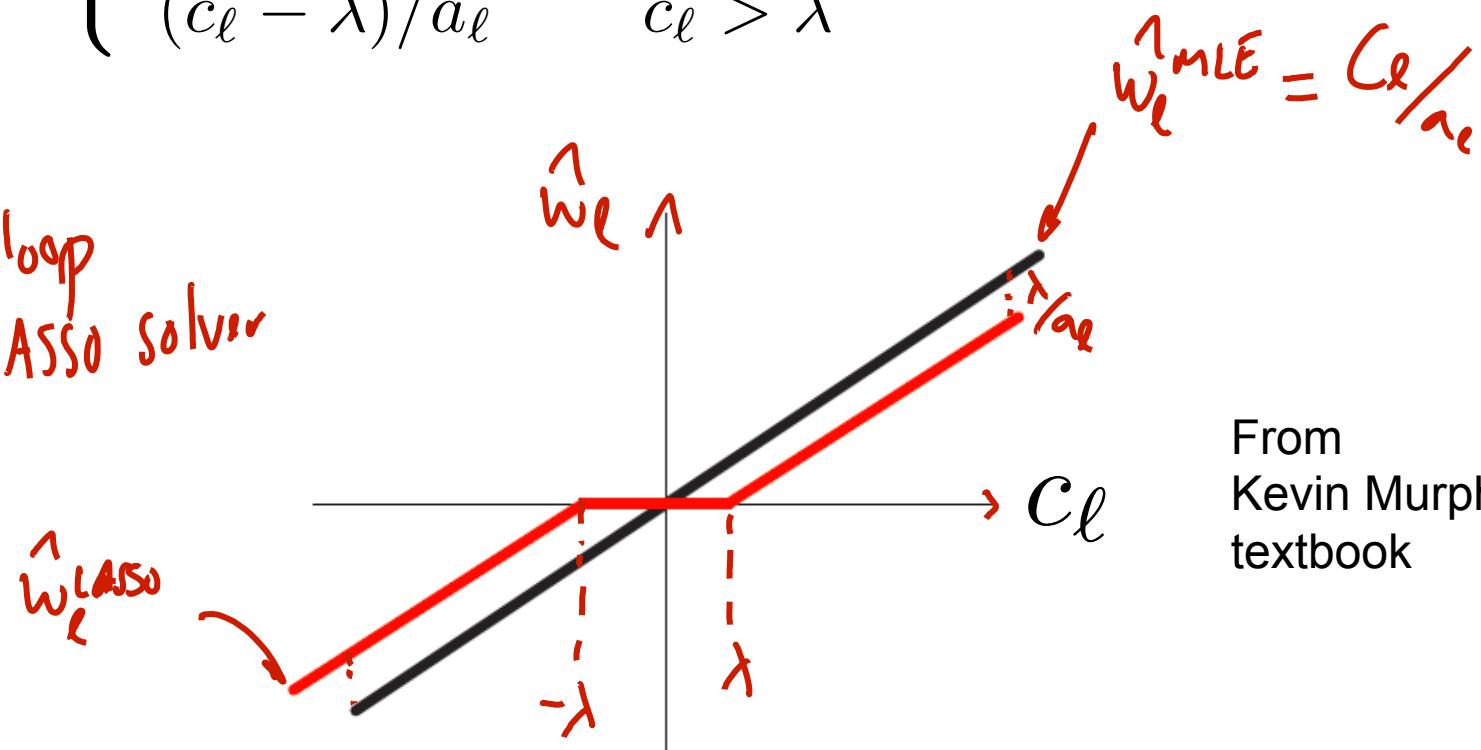
- Least angle regression and shrinkage, Efron et al. 2004

Soft Thresholding

$\lambda=0 \rightarrow \text{MLE solution is}$
 $\hat{w}_\ell^{\text{MLE}} = c_\ell/a_\ell$

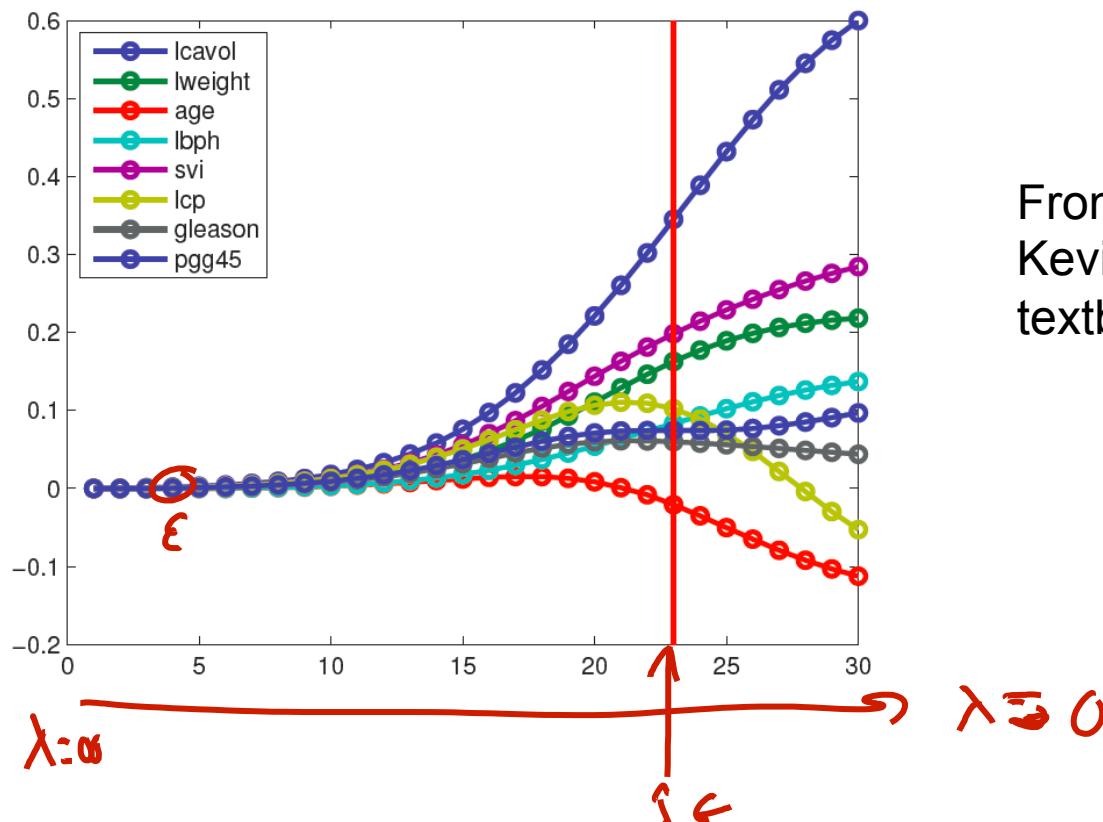
$$\hat{w}_\ell^{\text{LASSO}} = \begin{cases} (c_\ell + \lambda)/a_\ell & c_\ell < -\lambda \\ 0 & c_\ell \in [-\lambda, \lambda] \\ (c_\ell - \lambda)/a_\ell & c_\ell > \lambda \end{cases}$$

inner loop
of LASSO solver



From
Kevin Murphy
textbook

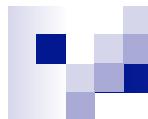
Recall: Ridge Coefficient Path



From
Kevin Murphy
textbook

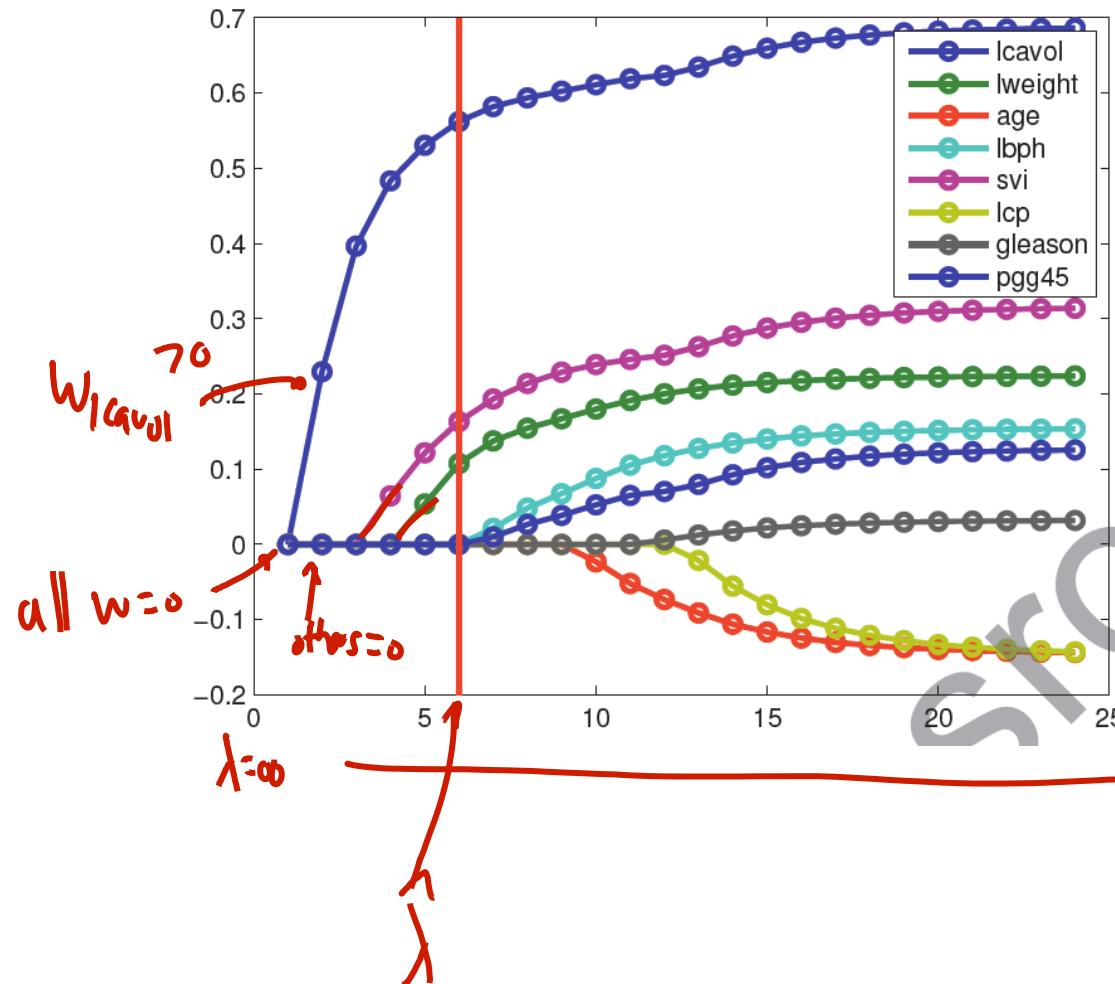
- Typical approach: select λ using cross validation

Now: LASSO Coefficient Path

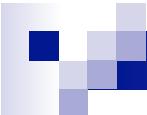


params can be
 $\emptyset \rightarrow \text{non zero} \rightarrow \emptyset$
as λ decreases

From
Kevin Murphy
textbook



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	

Wise

Wise

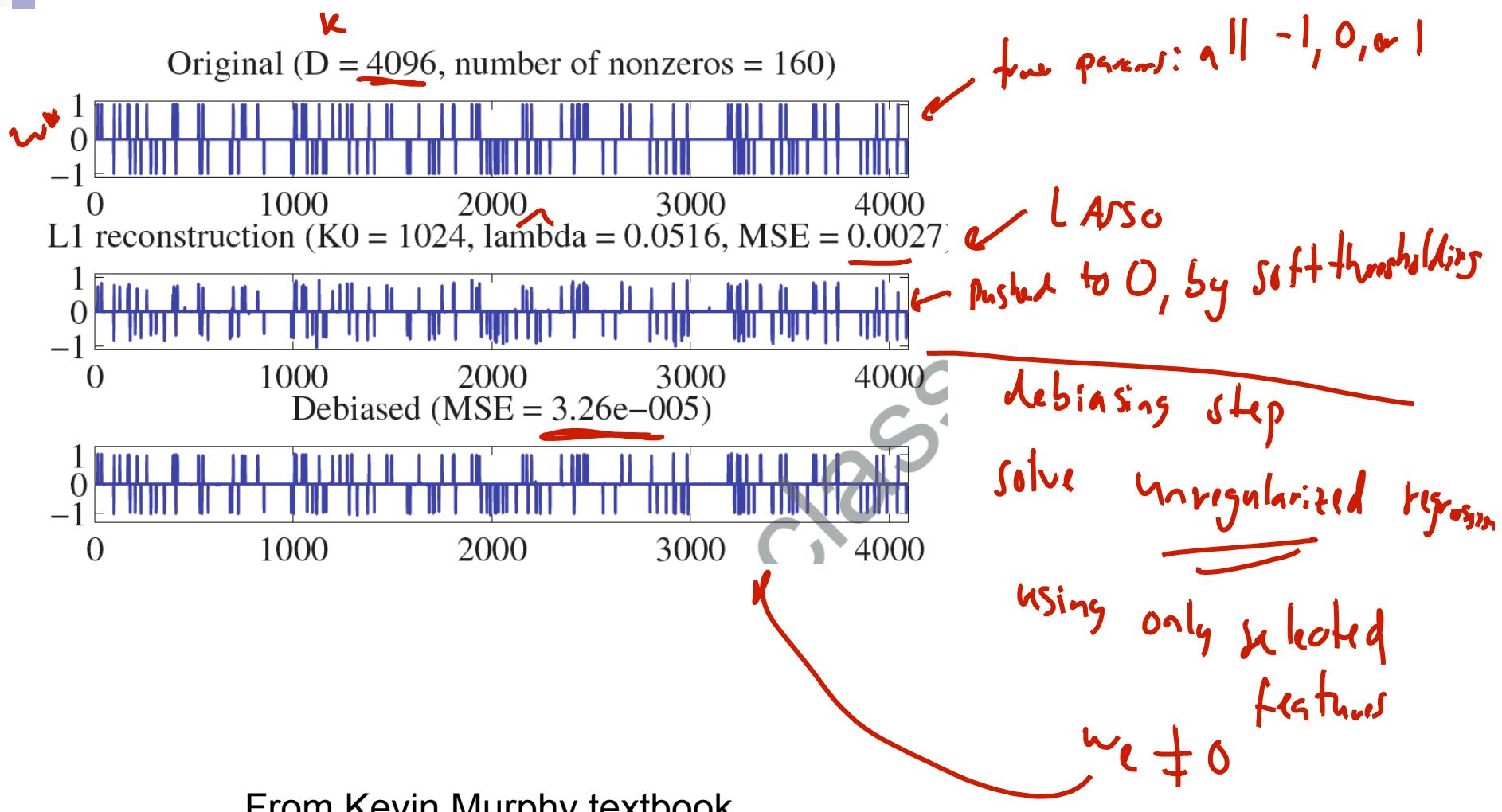
Wise

From
Rob
Tibshirani
slides



Debiasing

$$MSE = RSS$$



From Kevin Murphy textbook

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