

Training set error $w^* = \arg \min_w \sum_j \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$

- Given a dataset (Training data)
- Choose a loss function
 - e.g., squared error (L_2) for regression
- **Training set error:** For a particular set of parameters, loss function on training data:

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

Annotations:
 - \min_w : *min*
 - w : *hypothesis*
 - $\sum_{j=1}^{N_{train}}$: *sum train data*
 - $t(\mathbf{x}_j)$: *target*
 - $\sum_i w_i h_i(\mathbf{x}_j)$: *prediction*
 - 2 : *squared*

©2005-2013 Carlos Guestrin

3

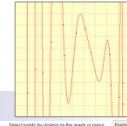
Training set error as a function of model complexity



©2005-2013 Carlos Guestrin

4

Prediction error



- Training set error can be poor measure of “quality” of solution
- **Prediction error:** We really care about error over all possible input points, not just training data:

want min

$$error_{true}(\mathbf{w}) = E_{\mathbf{x}} \left[\left(t(\mathbf{x}) - \sum_i w_i h_i(\mathbf{x}) \right)^2 \right]$$

all possible inputs

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

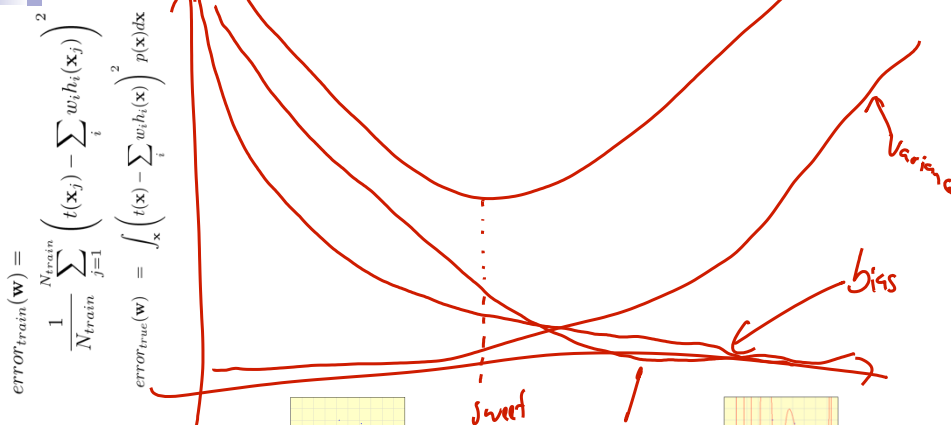
integrate/average

weighted by likelihood of \mathbf{x}

©2005-2013 Carlos Guestrin

5

Prediction error as a function of model complexity



©2005-2013 Carlos Guestrin

6

Computing prediction error

- Computing prediction

- Hard integral
- May not know $t(\mathbf{x})$ for every \mathbf{x}

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

- Monte Carlo integration (sampling approximation)

- Sample a set of i.i.d. points $\{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ from $p(\mathbf{x})$
- Approximate integral with sample average

$$error_{true}(\mathbf{w}) \approx \frac{1}{M} \sum_{j=1}^M \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

Sample approx of true $\mathbf{x}_j \sim p(\mathbf{x})$

©2005-2013 Carlos Guestrin

7

Why training set error doesn't approximate prediction error?

- Sampling approximation of prediction error:

$$error_{true}(\mathbf{w}) \approx \frac{1}{M} \sum_{j=1}^M \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

- Training error :

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

- Very similar equations!!!

- Why is training set a bad measure of prediction error???

©2005-2013 Carlos Guestrin

8

Why training set error doesn't approximate prediction error?

Because you cheated!!!

Training error good estimate for a single \mathbf{w} ,
But you optimized \mathbf{w} with respect to the training error,
and found \mathbf{w} that is good for this set of samples

**Training error is a (optimistically) biased
estimate of prediction error**

- Very similar equations!!!
 - Why is training set a bad measure of prediction error???

©2005-2013 Carlos Guestrin

9

Test set error

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

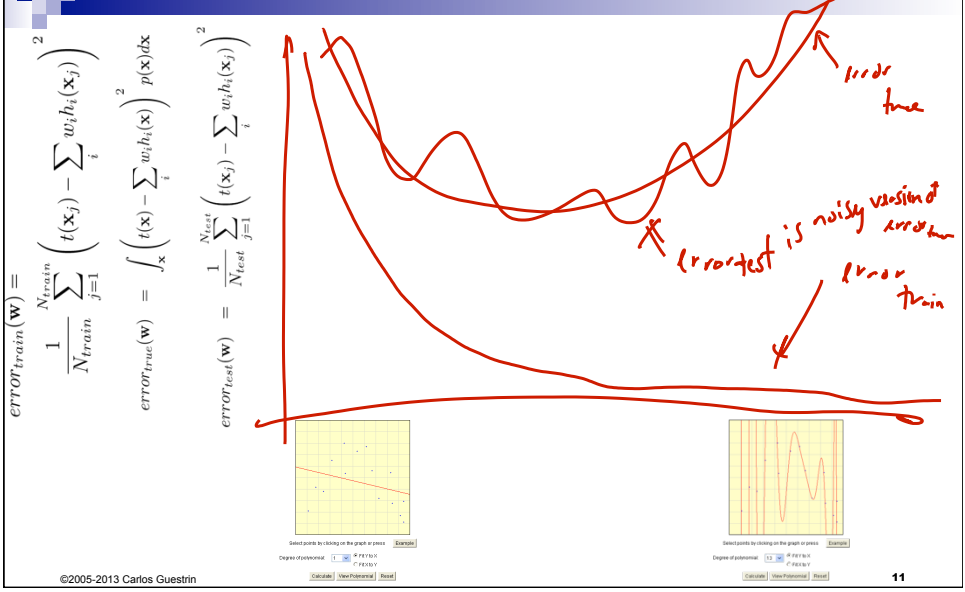
- Given a dataset, **randomly** split it into two parts:
 - Training data – $\{\mathbf{x}_1, \dots, \mathbf{x}_{N_{\text{train}}}\}$
 - Test data – $\{\mathbf{x}_1, \dots, \mathbf{x}_{N_{\text{test}}}\}$
- Use training data to optimize parameters \mathbf{w}
- **Test set error:** For the *final output* $\hat{\mathbf{w}}$, evaluate the error using:

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

©2005-2013 Carlos Guestrin

10

Test set error as a function of model complexity



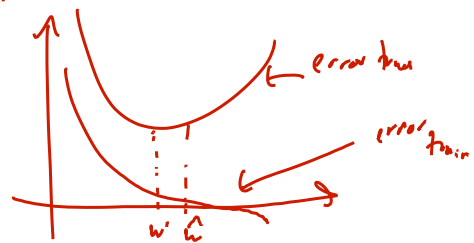
Overfitting

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

$$[error_{train}(\hat{\mathbf{w}}) < error_{train}(\mathbf{w}')] \wedge [error_{true}(\mathbf{w}') < error_{true}(\hat{\mathbf{w}})]$$

better on train data

worse in true error



How many points to I use for training/testing?

- Very hard question to answer!
 - Too few training points, learned $\hat{\mathbf{w}}$ is bad
 - Too few test points, you never know if you reached a good solution
- Bounds, such as Hoeffding's inequality can help:

$$P(|\hat{\theta} - \theta^*| \geq \epsilon) \leq 2e^{-2N\epsilon^2}$$

- More on this later this quarter, but still hard to answer
- Typically:
 - If you have a reasonable amount of data, pick test set "large enough" for a "reasonable" estimate of error, and use the rest for learning
 - If you have little data, then you need to pull out the big guns...
 - e.g., bootstrapping *(cross validation)*

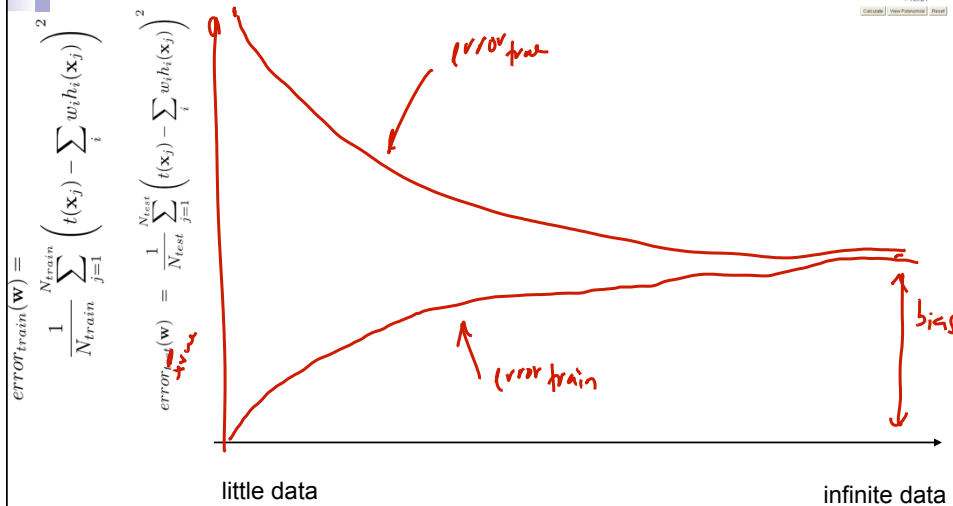
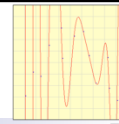
Error estimators

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

$$\text{error}_{\text{train}}(\mathbf{w}) = \frac{1}{N_{\text{train}}} \sum_{j=1}^{N_{\text{train}}} \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2 \quad \leftarrow \text{data for learning, optimistically biased solution}$$

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

Error as a function of number of training examples for a fixed model complexity



Error estimators

Be careful!!!

Test set only unbiased if you never never ever ever do any any any any learning on the test data

For example, if you use the test set to select the degree of the polynomial... no longer unbiased!!!
(We will address this problem later)

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

What you need to know

- True error, training error, test error

- Never learn on the test data
- Never learn on the test data
- Never learn on the test data
- Never learn on the test data
- Never learn on the test data



- Overfitting

Regularization

Machine Learning – CSE546

Carlos Guestrin

University of Washington

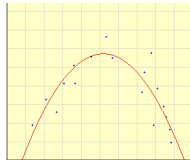
October 2, 2013

Regularization in Linear Regression

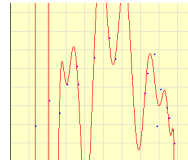
- Overfitting usually leads to very large parameter choices, e.g.:

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

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



less variance → smaller coeffs



large coeff

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

Ridge Regression

- Ameliorating issues with overfitting: *penalize large weights*

- New objective:

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

train error
||w_{1:k}||₂²

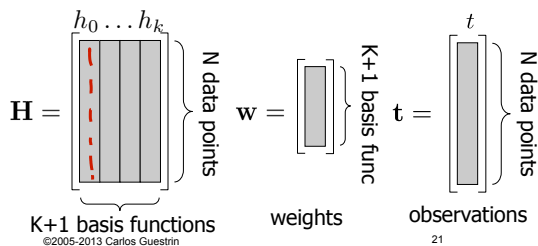
regularization
coeff

$$\lambda \geq 0$$

Ridge Regression in Matrix Notation

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

$$= \arg \min_{\mathbf{w}} \underbrace{(\mathbf{H}\mathbf{w} - \mathbf{t})^T (\mathbf{H}\mathbf{w} - \mathbf{t})}_{\text{residual error}} + \lambda \mathbf{w}^T \mathbf{I}_{0+k} \mathbf{w}$$



$$\mathbf{I}_{0+k} = \begin{pmatrix} 0 & & & & \\ & 1 & & & \\ & & \ddots & & \\ & & & 1 & \\ & & & & 0 & & \\ & & & & & \ddots & \\ & & & & & & 1 & \\ & & & & & & & \ddots \end{pmatrix}$$

Minimizing the Ridge Regression Objective

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

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

$$\nabla_{\mathbf{w}} F(\mathbf{w}) = 2\mathbf{H}^T (\mathbf{H}\mathbf{w} - \mathbf{t}) + 2\lambda \mathbf{I}_{0+k} \mathbf{w} = 0$$

$$\frac{d}{dw} \lambda w^2 = 2\lambda w$$

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

just matrix

adding $\lambda \mathbf{I}_{0+k}$ can only make $\mathbf{H}^T \mathbf{H}$ "more invertible" ops!!
 So regularization is also a numerical stabilization method

Shrinkage Properties

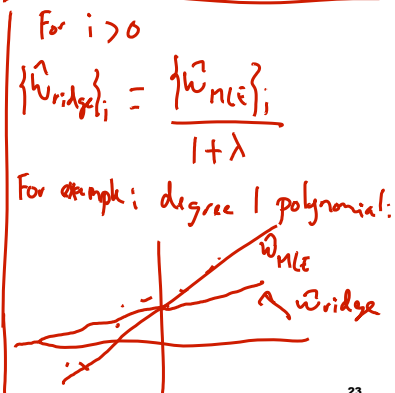
$$\hat{\mathbf{w}}_{ridge} = (H^T H + \lambda I_{0+k})^{-1} H^T \mathbf{t}$$

- If orthonormal features/basis: $H^T H = I$ $\hat{\mathbf{w}}_{MLE} = (H^T H)^{-1} H^T \mathbf{t}$

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

$$= \begin{pmatrix} 1 & & 0 \\ & 1+\lambda & \\ 0 & & 1+\lambda \dots \end{pmatrix}^{-1} H^T \mathbf{t}$$

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



©2005-2013 Carlos Guestrin

23

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 high regularization

- Smaller λ low regularization

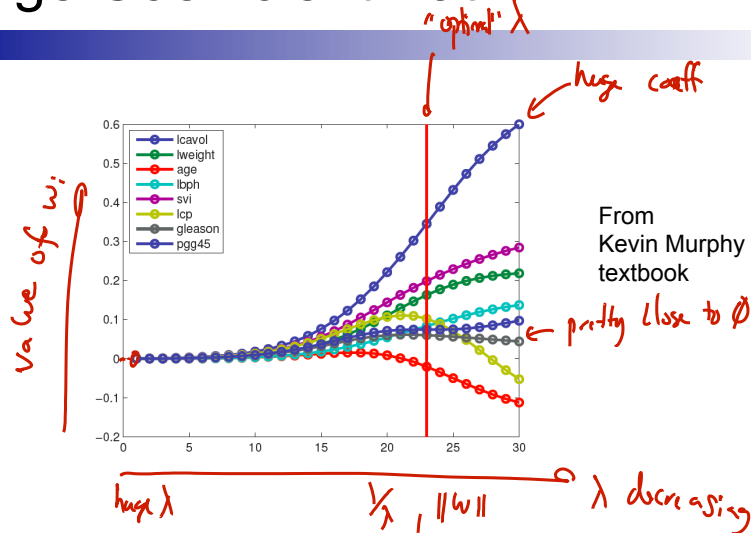
- As $\lambda \rightarrow 0$ $\hat{\mathbf{w}}_{ridge} \rightarrow \hat{\mathbf{w}}_{MLE}$

- As $\lambda \rightarrow \infty$ $\hat{\mathbf{w}}_{ridge} \rightarrow w_0 \leftarrow \text{average } t(x_j)$
flat line

©2005-2013 Carlos Guestrin

24

Ridge Coefficient Path

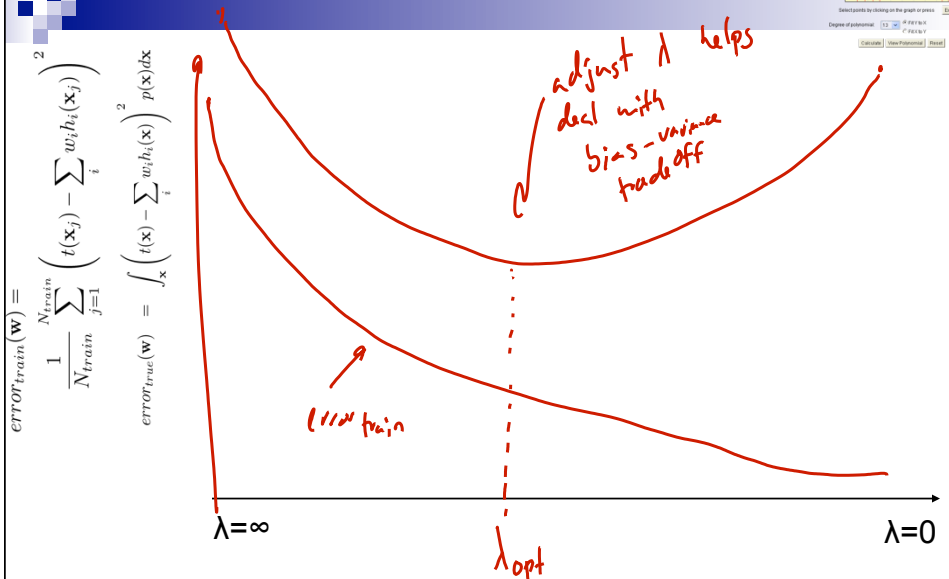


- Typical approach: select λ using cross validation, more on this later in the quarter

©2005-2013 Carlos Guestrin

25

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



©2005-2013 Carlos Guestrin

26

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

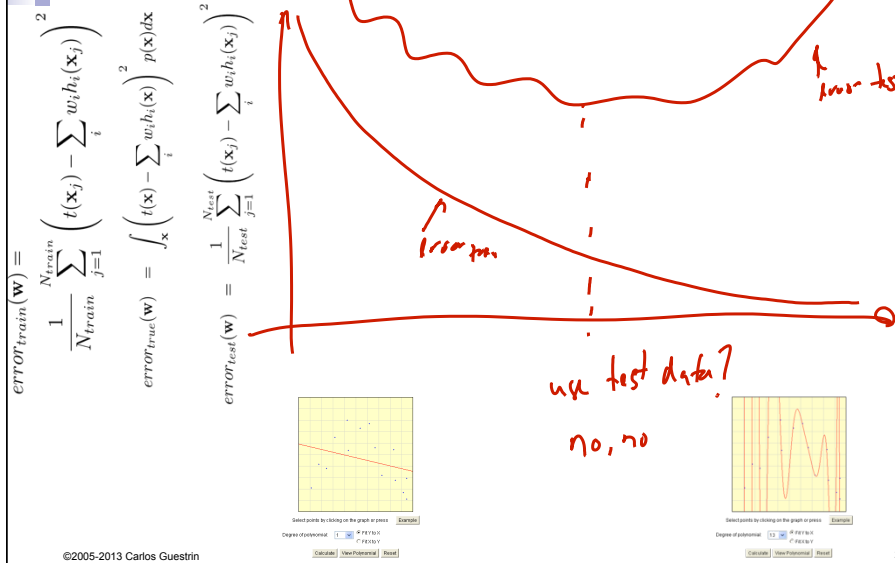
Machine Learning – CSE546

Carlos Guestrin

University of Washington

October 2, 2013

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. may not have enough data to separate some to pick λ

(LOO) Leave-one-out cross validation

- Consider a **validation set with 1 example**: *← choose λ on 1 training example*

- D – training data *← N data points*
- D_j – training data with j th data point moved to validation set

- Learn classifier h_{D_j} with D_j dataset** *← $N-1$ data points*

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

- Unbiased estimate of $\text{error}_{\text{true}}(h_{D_j})!$

$$= \mathbb{E}_x [(t - h_{D_j})^2]$$

$$\approx (t(\mathbf{x}_j) - h_{D_j}(\mathbf{x}_j))^2$$

- 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_j} *or regressor*

- Estimate error as:

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

©2005-2013 Carlos Guestrin

31

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!**

$$\text{error}_{\text{true}}(h_D) \approx \text{error}_{LOO}$$

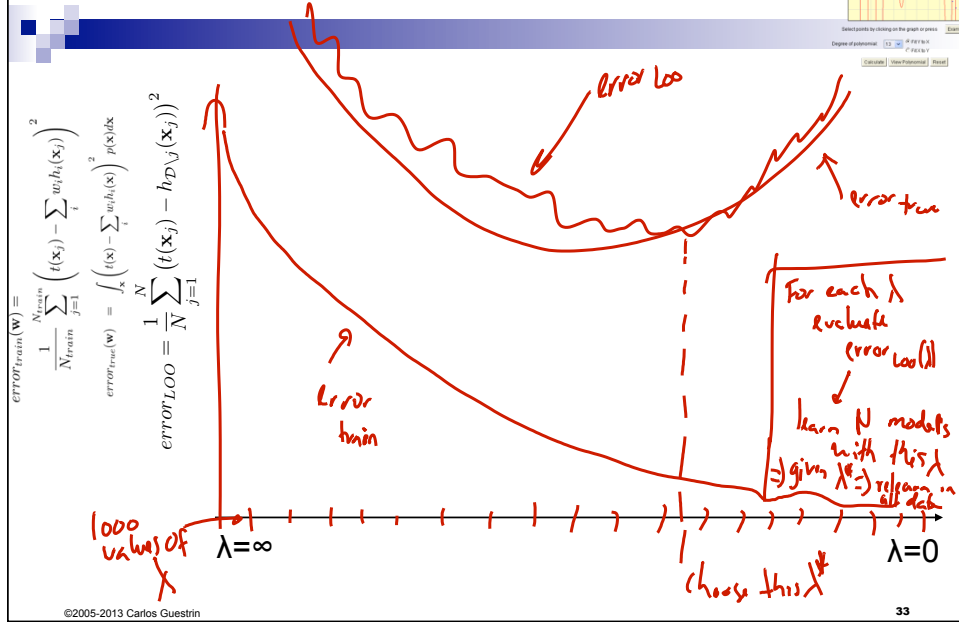
- Great news!**

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

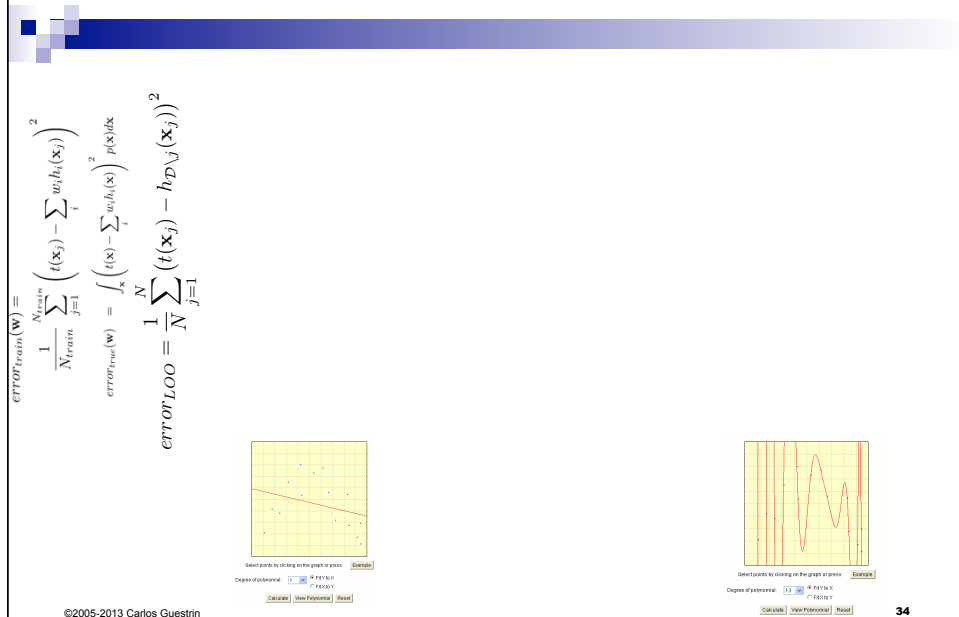
©2005-2013 Carlos Guestrin

32

Using LOO to Pick λ



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!!! of λ
 - If you have to do for each choice of basis functions, it will take fooooooreeve'!!!
- Solution 1: Preferred, but not usually possible
 - Find a cool trick to compute LOO (e.g., see homework)

©2005-2013 Carlos Guestrin

35

Solution 2 to complexity of computing LOO:

(More typical) **Use k -fold cross validation**

- Randomly divide training data into k equal parts LOO \equiv N -fold cross valid
 - D_1, \dots, D_k
- For each i regular or
 - Learn classifier h_{D_i} using data point not in D_i
 - Estimate error of h_{D_i} on validation set D_i :

$$error_{D_i} = \frac{1}{N} \sum_{\mathbf{x}_j \in D_i} (t(\mathbf{x}_j) - h_{D \setminus D_i}(\mathbf{x}_j))^2$$

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

$$error_{k\text{-fold}} = \frac{1}{k} \sum_{i=1}^k error_{D_i}$$

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

©2005-2013 Carlos Guestrin

36

