

Overfitting

Machine Learning – CSE546

Carlos Guestrin

University of Washington

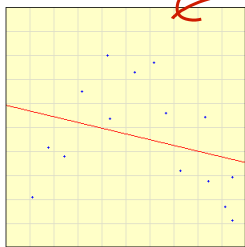
October 2, 2014

©2005-2014 Carlos Guestrin

1

Bias-Variance Tradeoff

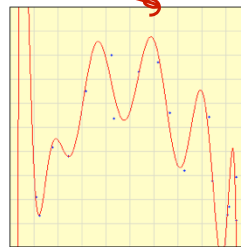
- Choice of hypothesis class introduces learning bias
 - More complex class → less bias
 - More complex class → more variance



Select points by clicking on the graph or press [Example](#)

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

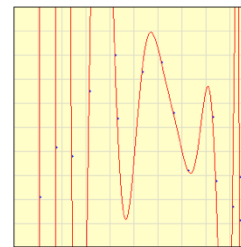
[Calculate](#) [View Polynomial](#) [Reset](#)



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](#)



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](#)

©2005-2014 Carlos Guestrin

2

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:

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

©2005-2014 Carlos Guestrin

3

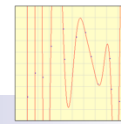
Training set error as a function of model complexity



©2005-2014 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:

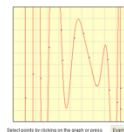
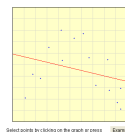
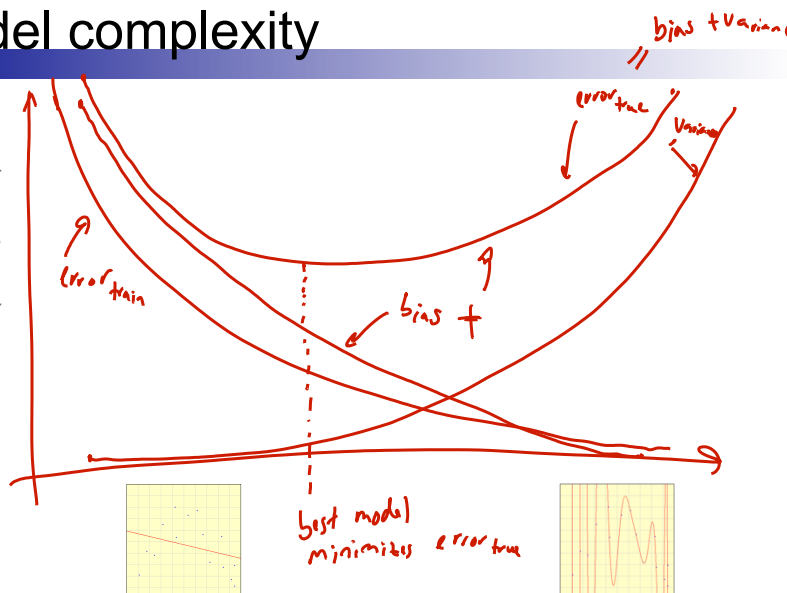
$$\begin{aligned}
 error_{true}(\mathbf{w}) &= E_{\mathbf{x}} \left[\left(t(\mathbf{x}) - \sum_i w_i h_i(\mathbf{x}) \right)^2 \right] \\
 &= \int_{\mathbf{x}} \left(t(\mathbf{x}) - \sum_i w_i h_i(\mathbf{x}) \right)^2 p(\mathbf{x}) d\mathbf{x}
 \end{aligned}$$

©2005-2014 Carlos Guestrin

5

Prediction error as a function of model complexity

$$\begin{aligned}
 error_{train}(\mathbf{w}) &= \frac{1}{N_{train}} \sum_{j=1}^{N_{train}} \left(t(\mathbf{x}_j) - \sum_k w_k h_k(\mathbf{x}_j) \right)^2 \\
 error_{true}(\mathbf{w}) &= \int_{\mathbf{x}} \left(t(\mathbf{x}) - \sum_k w_k h_k(\mathbf{x}) \right)^2 p(\mathbf{x}) d\mathbf{x}
 \end{aligned}$$



©2005-2014 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

Sample avg of error

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

©2005-2014 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???

©2005-2014 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???

©2005-2014 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:

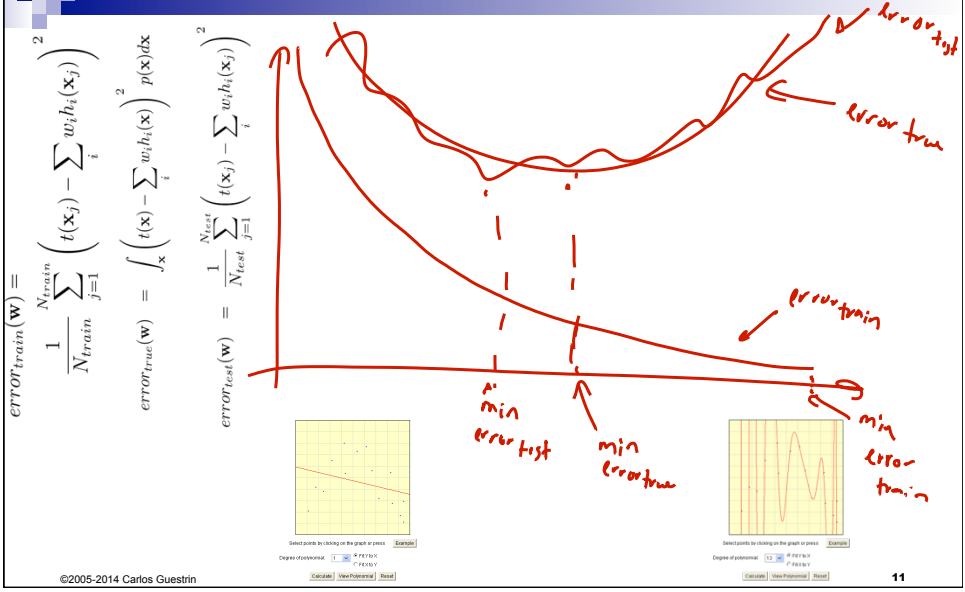
$$\underset{\substack{\text{error} \\ \text{error}_{\text{test}}(\hat{\mathbf{w}})}}{\text{error}_{\text{test}}(\hat{\mathbf{w}})} = \frac{1}{N_{\text{test}}} \sum_{j=1}^{N_{\text{test}}} \left(t(\mathbf{x}_j) - \sum_i \hat{w}_i h_i(\mathbf{x}_j) \right)^2$$

error_{true}(w)

©2005-2014 Carlos Guestrin

10

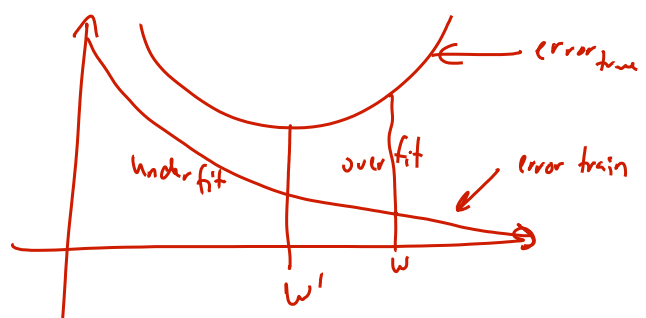
Test set error as a function of model complexity



Overfitting

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

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



How many points to I use for training/testing?

- Very hard question to answer!
 - Too few training points, learned \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

Error estimators

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

← gold standard

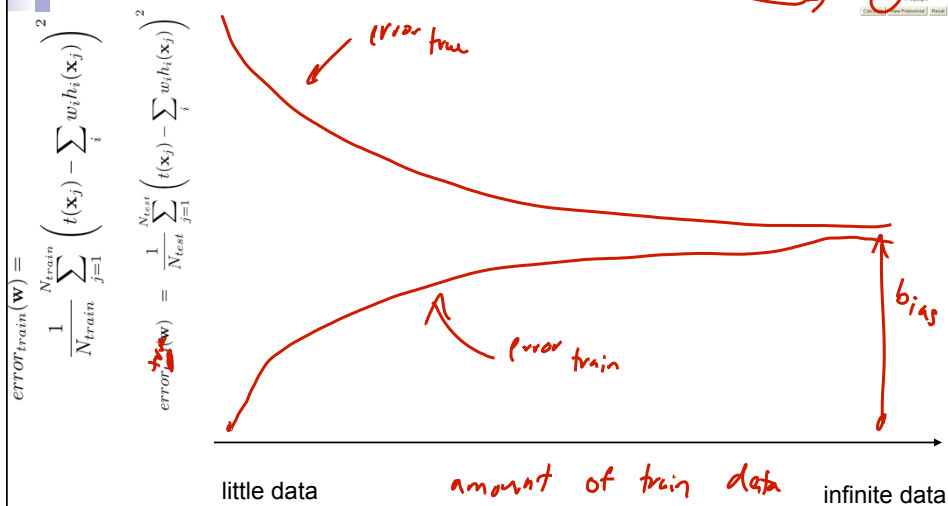
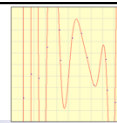
$$error_{train}(\mathbf{w}) = \frac{1}{N_{train}} \sum_{j=1}^{N_{train}} \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

← data for learning, optimistically biased

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

← Evaluate final model

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



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

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

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

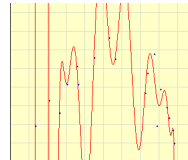
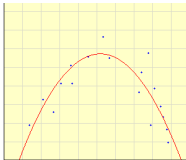
Machine Learning – CSE546
Carlos Guestrin
University of Washington
October 2, 2014

Regularization in Linear Regression

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

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

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



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

- “Shrinkage” method

→ smaller curvatures or smoother fits
 → a kind of bias we can believe in

Ridge Regression

- Ameliorating issues with overfitting: *penalizing large weights*

- New objective:

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

train error
regularization penalty

regularization coeff $\lambda \geq 0$
 regularization penalty
 practical note: don't regularize w_0

Shrinkage Properties

$$\hat{w}_{MLE} = (H^T H)^{-1} H^T t$$

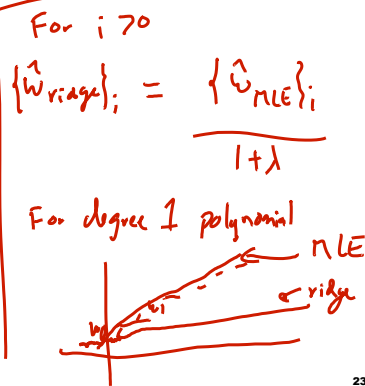
$$= \begin{pmatrix} 1 & 0 \\ 0 & I \end{pmatrix} H^T t$$

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

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

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

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



©2005-2014 Carlos Guestrin

23

Ridge Regression: Effect of Regularization

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

- Solution is indexed by the regularization parameter λ

- Larger λ higher regularization

- Smaller λ lower regularization

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

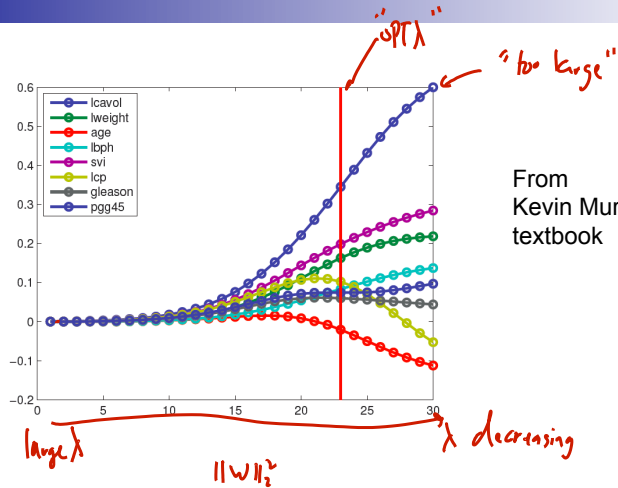
- As $\lambda \rightarrow \infty$ $\hat{w}_{ridge} \rightarrow w_0$ ← converges to avg of $t(x_j)$ in train data
flat line

©2005-2014 Carlos Guestrin

24

Ridge Coefficient Path

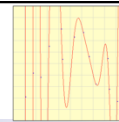
weights



From Kevin Murphy textbook

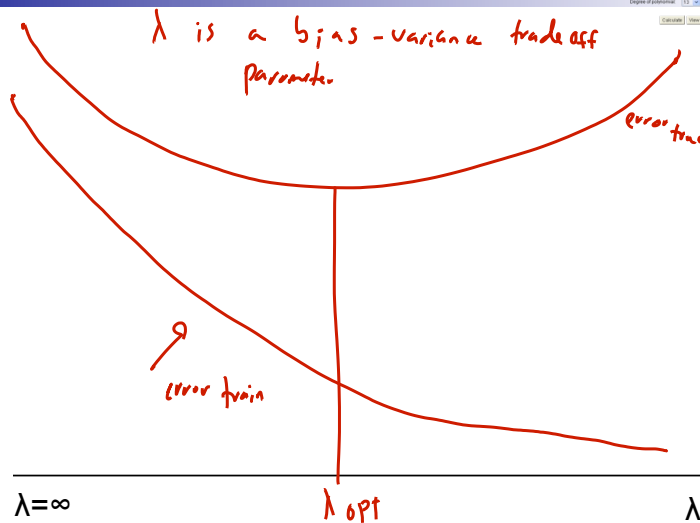
- Typical approach: select λ using cross validation, more on this later in the quarter

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



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

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



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

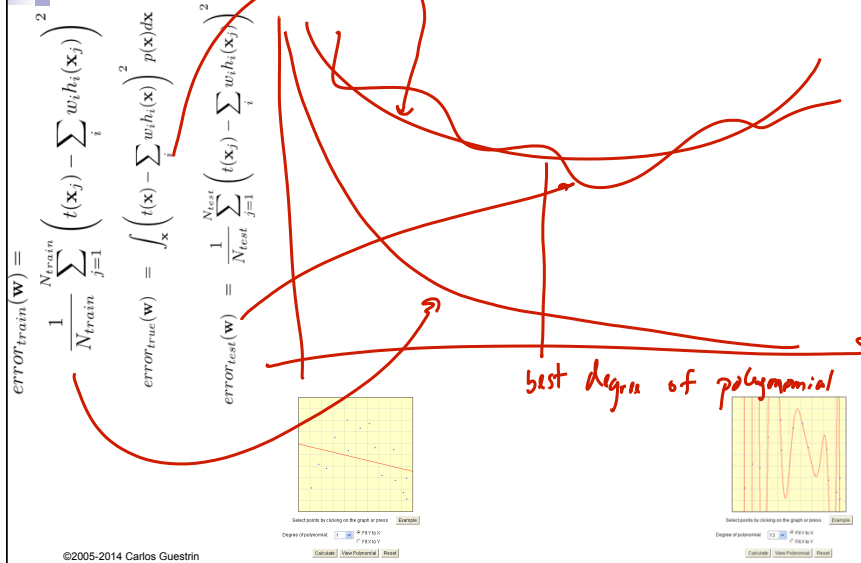
Machine Learning – CSE546

Carlos Guestrin

University of Washington

October 2, 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 learn on test data
 3. may not be enough data for a good train test split

(LOO) Leave-one-out cross validation

- Consider a **validation set with 1 example**: *← choose 1 using 1 train example*
 - D – training data *← N data points*
 - D_j – training data with j th data point moved to validation set *red square*
- **Learn classifier h_{D_j} with D_j dataset** *← using N-1 data points, for each j*
- **Estimate true error** as squared error on predicting $t(\mathbf{x}_j)$:
 - Unbiased estimate of error $\text{error}_{\text{true}}(h_{D_j})!$

$$E[(t(\mathbf{x}_j) - h_{D_j}(\mathbf{x}_j))^2] = \text{error}_{\text{true}}(h_{D_j})$$
 - 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}**
 - Estimate error as:

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



©2005-2014 Carlos Guestrin

31

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

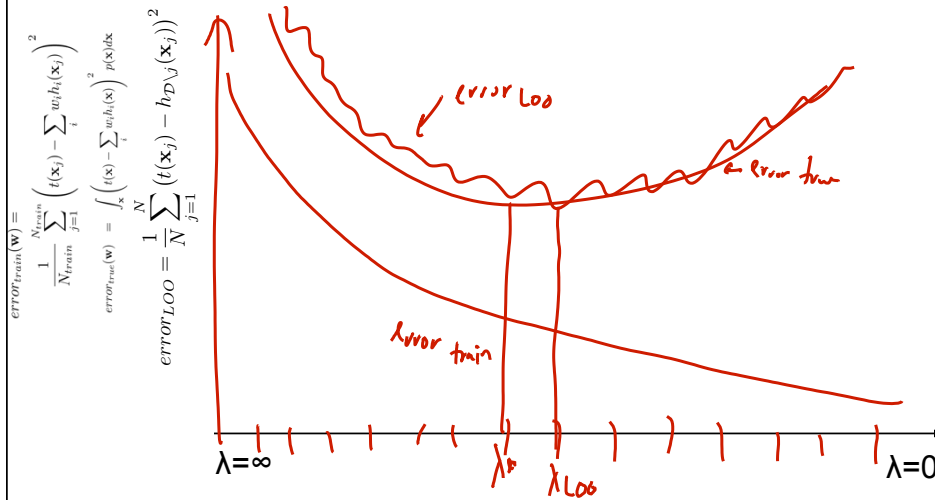
- When computing **LOOCV error, we only use $N-1$ data points**
 - So it's not estimate of true error of learning with N data points!
 - Usually pessimistic, though – learning with less data typically gives worse answer
- **LOO is almost unbiased!**

$$\text{error}_{\text{true}}(h_D) \approx \text{error}_{\text{LOO}}$$
- **Great news!**
 - **Use LOO error for model selection!!!**
 - **E.g., picking λ**

©2005-2014 Carlos Guestrin

32

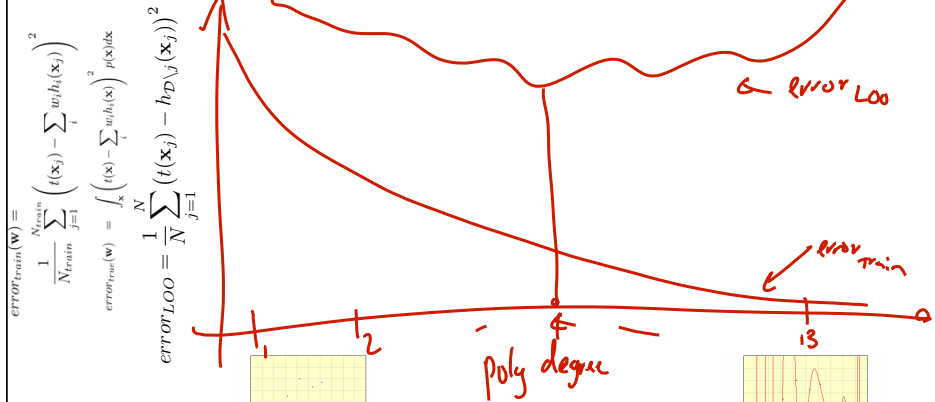
Using LOO to Pick λ



©2005-2014 Carlos Guestrin

33

Using LOO error for model selection



©2005-2014 Carlos Guestrin

34

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!!! or 1
 - 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)

©2005-2014 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
= N-fold
 - 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 :

$$error_{D_i} = \frac{1}{N} \sum_{\mathbf{x}_j \in D_i} (t(\mathbf{x}_j) - h_{D \setminus D_i}(\mathbf{x}_j))^2$$
- 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, $k = 10$ ☺

©2005-2014 Carlos Guestrin

36

