

Overfitting

Machine Learning – CSE446

Carlos Guestrin

University of Washington

April 8, 2013

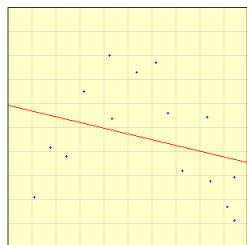
©2005-2013 Carlos Guestrin

1

Bias-Variance Tradeoff

- Choice of hypothesis class introduces learning bias

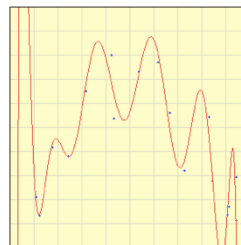
- ☐ More complex class \rightarrow less bias
- ☐ More complex class \rightarrow more variance



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

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

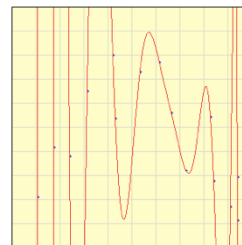
[Calculate](#) [View Polynomial](#) [Reset](#)



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](#)

©2005-2013 Carlos Guestrin

2

Training 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 (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:

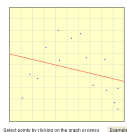
$$error_{train}(\mathbf{w}) = \frac{1}{N_{train}} \sum_{j=1}^{N_{train}} \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

©2005-2013 Carlos Guestrin

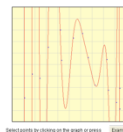
3

Training set error as a function of model complexity

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



©2005-2013 Carlos Guestrin



©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:

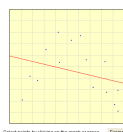
$$\begin{aligned} error_{true}(\mathbf{w}) &= E_{\mathbf{x}} \left[\left(t(\mathbf{x}) - \sum_i w_i h_i(\mathbf{x}) \right)^2 \right] \\ &= \int_{\mathbf{x}} \left(t(\mathbf{x}) - \sum_i w_i h_i(\mathbf{x}) \right)^2 p(\mathbf{x}) d\mathbf{x} \end{aligned}$$

©2005-2013 Carlos Guestrin

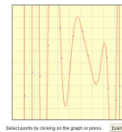
5

Prediction error as a function of model complexity: Bias/Variance tradeoff

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



Select points to clicking on the graph or press [Random](#)
Degree of polynomial: [Fit](#) [Fit to 0](#) [Fit to 1](#)
[Calculate](#) [View Polynomial](#) [Reset](#)



Select points to clicking on the graph or press [Random](#)
Degree of polynomial: [Fit](#) [Fit to 0](#) [Fit to 1](#)
[Calculate](#) [View Polynomial](#) [Reset](#)

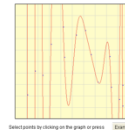
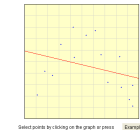
©2005-2013 Carlos Guestrin

6

Prediction error as a function of model complexity: train v. true error

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

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



©2005-2013 Carlos Guestrin

7

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

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

©2005-2013 Carlos Guestrin

8

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

9

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

10

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

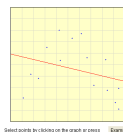
11

Test set error as a function of model complexity

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

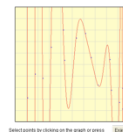
$$error_{true}(\mathbf{w}) = \int_{\mathbf{x}} \left(t(\mathbf{x}) - \sum_i w_i h_i(\mathbf{x}) \right)^2 p(\mathbf{x}) d\mathbf{x}$$

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



Graph of performance: ☐ Error ☒ R-squared

[Calculate](#) [Reset](#) [Zoom](#)



Graph of performance: ☒ Error ☐ R-squared

[Calculate](#) [Reset](#) [Zoom](#)

©2005-2013 Carlos Guestrin

12

Overfitting

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

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

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

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

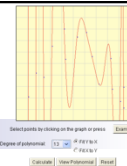
$$error_{train}(\mathbf{w}) = \frac{1}{N_{train}} \sum_{j=1}^{N_{train}} \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

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

©2005-2013 Carlos Guestrin

15

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



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

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

little data

infinite data

©2005-2013 Carlos Guestrin

16

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 in the quarter)

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

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 – CSE446

Carlos Guestrin

University of Washington

April 8, 2013

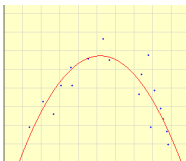
©2005-2013 Carlos Guestrin

19

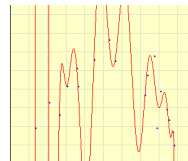
Regularization in Linear Regression

- 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

©2005-2013 Carlos Guestrin

20

Ridge Regression

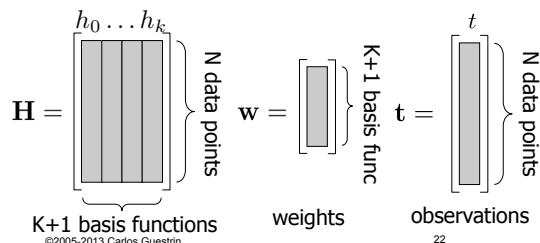
- Ameliorating issues with overfitting:
- New objective:

©2005-2013 Carlos Guestrin

21

Ridge Regression in Matrix Notation

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



©2005-2013 Carlos Guestrin

22

Minimizing the Ridge Regression Objective

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

©2005-2013 Carlos Guestrin

23

Shrinkage Properties

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

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

©2005-2013 Carlos Guestrin

24

Ridge Regression: Effect of Regularization

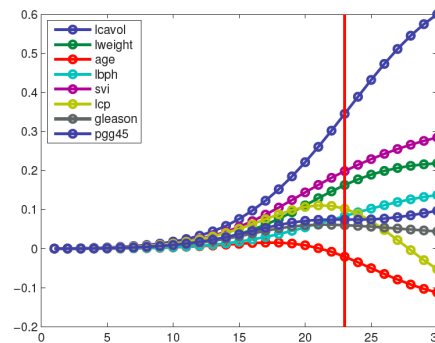
$$\hat{\mathbf{w}}_{ridge} = \arg \min_w \sum_{j=1}^N \left(t(x_j) - (w_0 + \sum_{i=1}^k w_i h_i(x_j)) \right)^2 + \lambda \sum_{i=1}^k w_i^2$$

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

©2005-2013 Carlos Guestrin

25

Ridge Coefficient Path



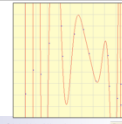
From
Kevin Murphy
textbook

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

©2005-2013 Carlos Guestrin

26

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



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

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

$\lambda = \infty$

$\lambda = 0$

©2005-2013 Carlos Guestrin

27

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

©2005-2013 Carlos Guestrin

28