



# Linear Regression

Machine Learning – CSE446

Carlos Guestrin


University of Washington

April 5, 2013

©2005-2013 Carlos Guestrin

1

## Prediction of continuous variables

- 
- Billionaire sayz: Wait, that's not what I meant!
  - You sayz: Chill out, dude.
  - He sayz: I want to predict a continuous variable for continuous inputs: I want to predict salaries from GPA.
  - You sayz: **I can regress that...**

©2005-2013 Carlos Guestrin

2

# The regression problem

- **Instances:**  $\langle \mathbf{x}_j, t_j \rangle$
- **Learn:** Mapping from  $\mathbf{x}$  to  $t(\mathbf{x})$
- **Hypothesis space:**  $H = \{h_1, \dots, h_K\}$ 
  - Given, basis functions
  - Find coeffs  $\mathbf{w} = \{w_1, \dots, w_K\}$   $\underbrace{t(\mathbf{x})}_{\text{data}} \approx \hat{f}(\mathbf{x}) = \sum_i w_i h_i(\mathbf{x})$
  - Why is this called linear regression???
  - model is linear in the parameters

- Precisely, minimize the **residual squared error**:

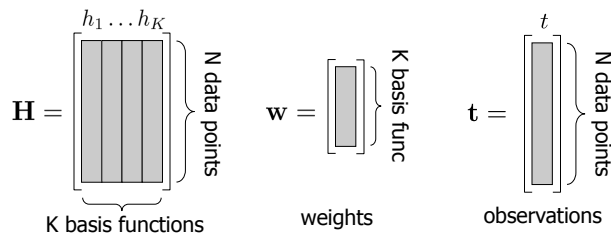
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_j \left( t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

©2005-2013 Carlos Guestrin

## The regression problem in matrix notation

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

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



©2005-2013 Carlos Guestrin

# Minimizing the Residual

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

©2005-2013 Carlos Guestrin

5

## Regression solution = simple matrix operations

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

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

$$\text{where } \mathbf{A} = \mathbf{H}^T \mathbf{H} = \underbrace{\begin{bmatrix} \square & \square & \square \\ \square & \square & \square \\ \square & \square & \square \end{bmatrix}}_{\substack{k \times k \text{ matrix} \\ \text{for } k \text{ basis functions}}} \quad \mathbf{b} = \mathbf{H}^T \mathbf{t} = \underbrace{\begin{bmatrix} \square \\ \square \\ \square \end{bmatrix}}_{k \times 1 \text{ vector}}$$

©2005-2013 Carlos Guestrin

6

## But, why?

- Billionaire (again) says: Why sum squared error???
- You say: Gaussians, Dr. Gateson, Gaussians...
- Model: prediction is linear function plus Gaussian noise
  - $t(\mathbf{x}) = \sum_i w_i h_i(\mathbf{x}) + \varepsilon_{\mathbf{x}}$

- Learn  $\mathbf{w}$  using MLE

$$P(t \mid \mathbf{x}, \mathbf{w}, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{[t - \sum_i w_i h_i(\mathbf{x})]^2}{2\sigma^2}}$$

©2005-2013 Carlos Guestrin

7

## Maximizing log-likelihood

**Maximize:**

$$\ln P(\mathcal{D} \mid \mathbf{w}, \sigma) = \ln \left( \frac{1}{\sigma\sqrt{2\pi}} \right)^N \prod_{j=1}^N e^{-\frac{[t_j - \sum_i w_i h_i(\mathbf{x}_j)]^2}{2\sigma^2}}$$

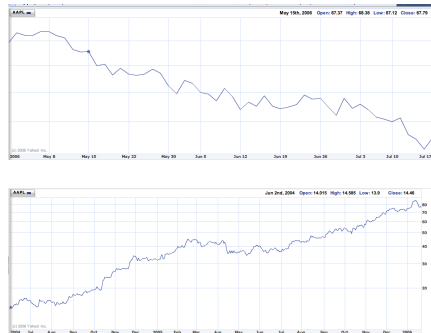
**Least-squares Linear Regression is MLE for Gaussians!!!**

©2005-2013 Carlos Guestrin

8

## Applications Corner 1

- Predict stock value over time from
  - past values
  - other relevant vars
    - e.g., weather, demands, etc.

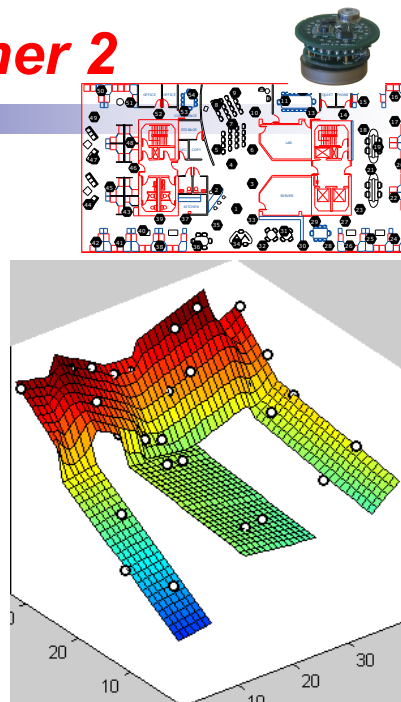


©2005-2013 Carlos Guestrin

9

## Applications Corner 2

- Measure temperatures at some locations
- Predict temperatures throughout the environment



[Guestrin et al. '04]

©2005-2013 Carlos Guestrin

10

## ***Applications Corner 3***

- Predict when a sensor will fail
  - based several variables
    - age, chemical exposure, number of hours used,...

©2005-2013 Carlos Guestrin

11

## **Bias-Variance Tradeoff**

Machine Learning – CSE446  
Carlos Guestrin  
University of Washington  
April 5, 2013

©2005-2013 Carlos Guestrin

12

## Bias-Variance tradeoff – Intuition

- Model too “simple” → does not fit the data well
  - A biased solution
- Model too complex → small changes to the data, solution changes a lot
  - A high-variance solution

©2005-2013 Carlos Guestrin

13

## (Squared) Bias of learner

- Given dataset  $D$  with  $N$  samples, learn function  $h_D(x)$
- If you sample a different dataset  $D'$  with  $N$  samples, you will learn different  $h_{D'}(x)$
- **Expected hypothesis:**  $E_D[h_D(x)]$
- **Bias:** difference between what you expect to learn and truth
  - Measures how well you expect to represent true solution
  - Decreases with more complex model
  - Bias<sup>2</sup> at one point  $x$ :
  - Average Bias<sup>2</sup>:

©2005-2013 Carlos Guestrin

14

# Variance of learner

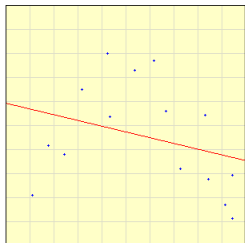
- Given dataset  $D$  with  $N$  samples, learn function  $h_D(x)$
- If you sample a different dataset  $D'$  with  $N$  samples, you will learn different  $h_{D'}(x)$
- **Variance:** difference between what you expect to learn and what you learn from a particular dataset
  - Measures how sensitive learner is to specific dataset
  - Decreases with simpler model
  - Variance at one point  $x$ :
  - Average variance:

©2005-2013 Carlos Guestrin

15

# Bias-Variance Tradeoff

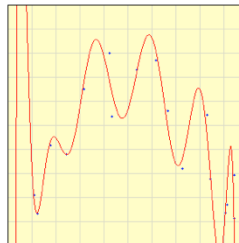
- Choice of hypothesis class introduces learning bias
  - More complex class  $\rightarrow$  less bias
  - More complex class  $\rightarrow$  more variance



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

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

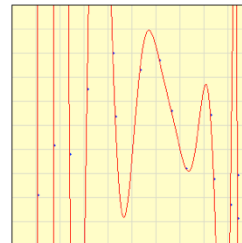
[Calculate](#) [View Polynomial](#) [Reset](#)



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

16



## Bias-Variance Decomposition of Error

$$\bar{h}_N(x) = E_D[h_D(x)]$$

- Expected mean squared error:  $\text{MSE} = E_D \left[ E_x \left[ (t(x) - h_D(x))^2 \right] \right]$
- To simplify derivation, drop  $x$ :
- Expanding the square:

## Moral of the Story: Bias-Variance Tradeoff Key in ML

- Error can be decomposed:
$$\begin{aligned}\text{MSE} &= E_D \left[ E_x \left[ (t(x) - h_D(x))^2 \right] \right] \\ &= E_x \left[ (t(x) - \bar{h}_N(x))^2 \right] + E_D \left[ E_x \left[ (\bar{h}_N(x) - h_D(x))^2 \right] \right]\end{aligned}$$
- Choice of hypothesis class introduces learning bias
  - More complex class  $\rightarrow$  less bias
  - More complex class  $\rightarrow$  more variance

# What you need to know

- Regression

- ☐ Basis function = features
- ☐ Optimizing sum squared error
- ☐ Relationship between regression and Gaussians

- Bias-variance trade-off

- Play with Applet

©2005-2013 Carlos Guestrin

19

## Overfitting

Machine Learning – CSE446

Carlos Guestrin

University of Washington

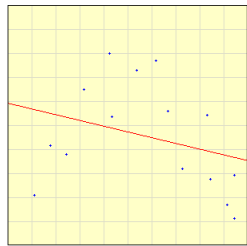
April 5, 2013

©2005-2013 Carlos Guestrin

20

# Bias-Variance Tradeoff

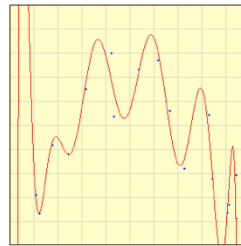
- Choice of hypothesis class introduces learning bias
  - More complex class → less bias
  - More complex class → more variance



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

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

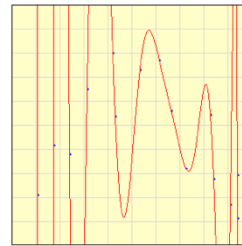
[Calculate](#) [View Polynomial](#) [Reset](#)



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](#)

©2005-2013 Carlos Guestrin

21

# 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:

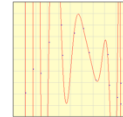
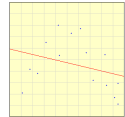
$$error_{train}(\mathbf{w}) = \frac{1}{N_{train}} \sum_{j=1}^{N_{train}} \left( t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

©2005-2013 Carlos Guestrin

22

# 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

23

## 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:

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

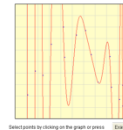
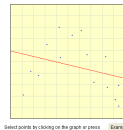
©2005-2013 Carlos Guestrin

24

# Prediction error as a function of model complexity

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

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



©2005-2013 Carlos Guestrin

25

## 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

26

## 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

27

## 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

28

## 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

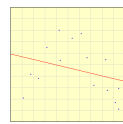
29

## Test set error as a function of model complexity

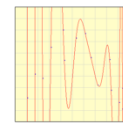
$$error_{train}(\mathbf{w}) = \frac{1}{N_{train}} \sum_{j=1}^{N_{train}} \left( t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

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

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



Graph of polynomial  $f(x) = 0.5x^2 + 1.5x + 0.5$   
[Zoom](#) [Reset](#) [Help](#)



Graph of polynomial  $f(x) = 0.5x^2 + 1.5x + 0.5$   
[Zoom](#) [Reset](#) [Help](#)

©2005-2013 Carlos Guestrin

30

# 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

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

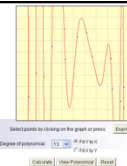
$$error_{train}(\mathbf{w}) = \frac{1}{N_{train}} \sum_{j=1}^{N_{train}} \left( t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

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

©2005-2013 Carlos Guestrin

33

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



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

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

little data

infinite data

©2005-2013 Carlos Guestrin

34

# 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