

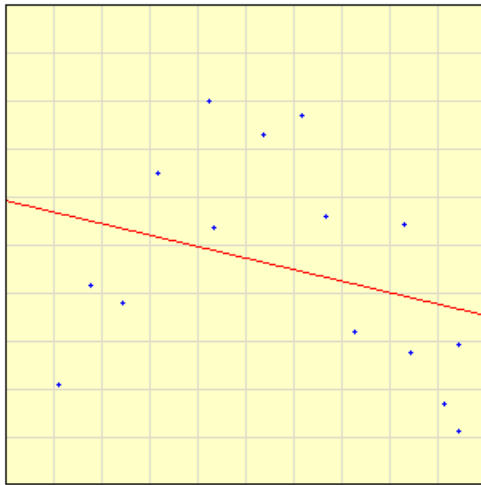
# CSE 575: Statistical Machine Learning

Jingrui He  
CIDSE, ASU

# Overfitting

# Bias-Variance Tradeoff

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

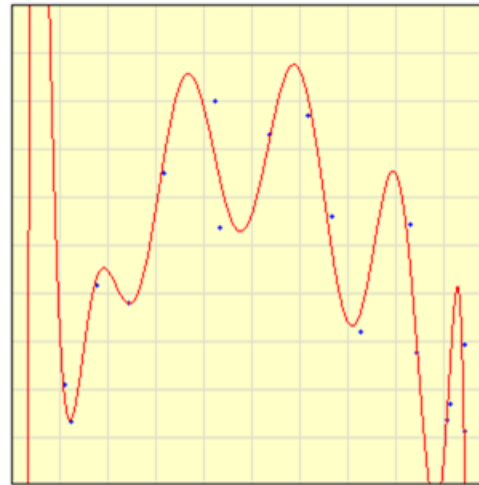


Select points by clicking on the graph or press

Example

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

Calculate View Polynomial Reset

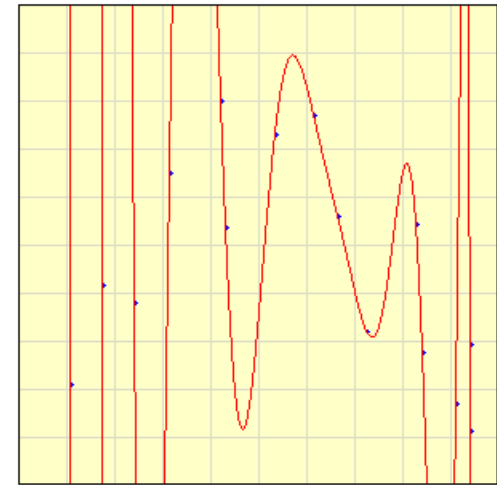


Select points by clicking on the graph or press

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

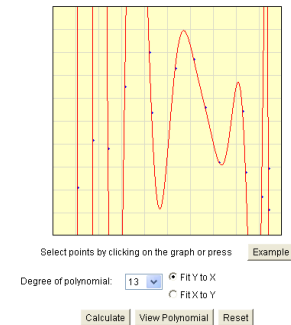
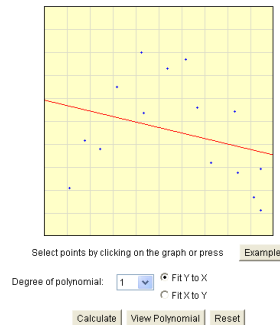
# Training set error $\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_j \left( t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$

- 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:

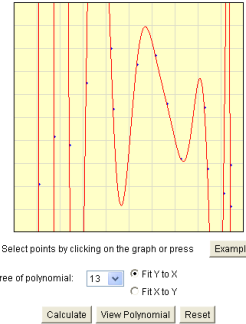
$$error_{train}(\mathbf{w}) = \frac{1}{N_{train}} \sum_{j=1}^{N_{train}} \left( t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

# Training set error as a function of model complexity

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



# Prediction error

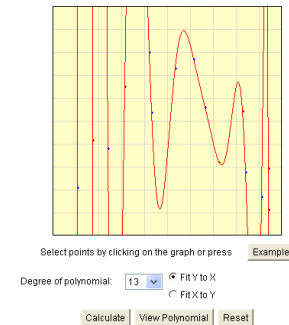
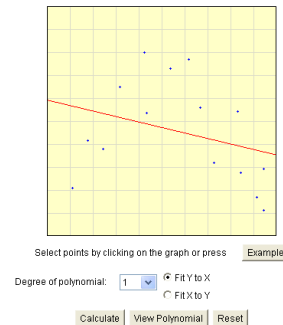


- Training set error can be poor measure of “quality” of solution
- **Prediction error:** We really care about error over all possible input points, not just training data:

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

# Prediction error as a function of model complexity

$$\begin{aligned} error_{train}(\mathbf{w}) &= \frac{1}{N_{train}} \sum_{j=1}^{N_{train}} \left( t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2 \\ error_{truc}(\mathbf{w}) &= \int_{\mathbf{x}} \left( t(\mathbf{x}) - \sum_i w_i h_i(\mathbf{x}) \right)^2 p(\mathbf{x}) d\mathbf{x} \end{aligned}$$



# 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

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



# 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???

# 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???

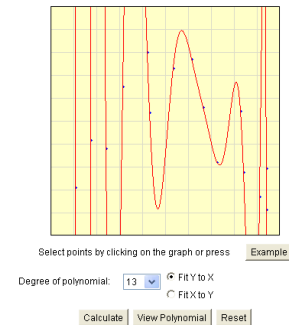
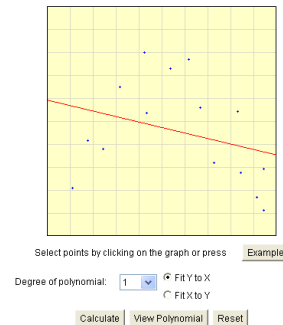
# Test set error $\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_j \left( t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$

- 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 solution***  $\mathbf{w}^*$ , evaluate the error using:

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

# Test set error as a function of model complexity

$$\begin{aligned} error_{train}(\mathbf{w}) &= \frac{1}{N_{train}} \sum_{j=1}^{N_{train}} \left( t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2 \\ error_{truc}(\mathbf{w}) &= \int_{\mathbf{x}} \left( t(\mathbf{x}) - \sum_i w_i h_i(\mathbf{x}) \right)^2 p(\mathbf{x}) d\mathbf{x} \\ error_{test}(\mathbf{w}) &= \frac{1}{N_{test}} \sum_{j=1}^{N_{test}} \left( t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2 \end{aligned}$$



# 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:

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

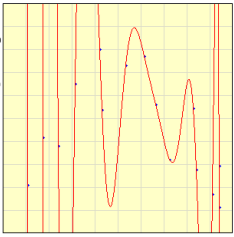
# Error estimators

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

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

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

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



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](#)

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

little data

infinite data

# Error estimators

**Be careful!!!**

Test set only unbiased if you never never never never  
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 semester)

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



# What's (supervised) learning more formally

- Given:
  - **Dataset:** Instances  $\{\langle \mathbf{x}_1; t(\mathbf{x}_1) \rangle, \dots, \langle \mathbf{x}_N; t(\mathbf{x}_N) \rangle\}$ 
    - e.g.,  $\langle \mathbf{x}_i; t(\mathbf{x}_i) \rangle = \langle (\text{GPA}=3.9, \text{IQ}=120, \text{MLscore}=99); 150\text{K} \rangle$
  - **Hypothesis space:**  $H$ 
    - e.g., polynomials of degree 8
  - **Loss function:** measures quality of hypothesis  $h \in H$ 
    - e.g., squared error for regression
- Obtain:
  - **Learning algorithm:** obtain  $h \in H$  that minimizes loss function
    - e.g., using matrix operations for regression
    - Want to minimize prediction error, but can only minimize error in dataset

# Types of (supervised) learning problems

- **Regression**, e.g.,
  - **dataset**:  $\langle \text{position; temperature} \rangle$
  - **hypothesis space**:
  - **Loss function**:
- **Density estimation**, e.g.,
  - **dataset**:  $\langle \text{grades} \rangle$
  - **hypothesis space**:
  - **Loss function**:
- **Classification**, e.g.,
  - **dataset**:  $\langle \text{brain image; \{verb v. noun\}} \rangle$
  - **hypothesis space**:
  - **Loss function**:

# Learning is (simply) function approximation!

- The general (supervised) learning problem:
  - Given some data (including features), hypothesis space, loss function
  - Learning is no magic!
  - Simply trying to find a function that fits the data
- **Regression**
- **Density estimation**
- **Classification**
- (Not surprisingly) Seemly different problem, very similar solutions...