

What's learning, revisited

Overfitting

Generative versus Discriminative

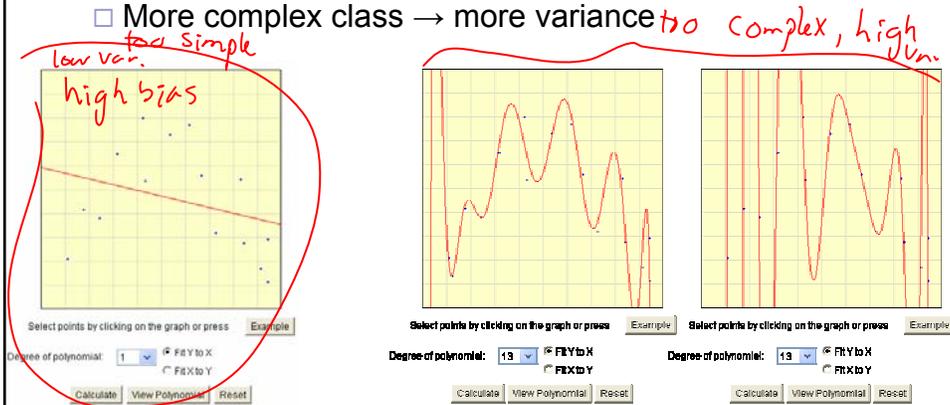
Logistic Regression

Machine Learning – 10701/15781
Carlos Guestrin
Carnegie Mellon University
September 19th, 2007

©Carlos Guestrin 2005-2007

Bias-Variance Tradeoff

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



©Carlos Guestrin 2005-2007

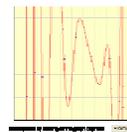
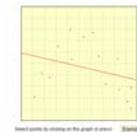
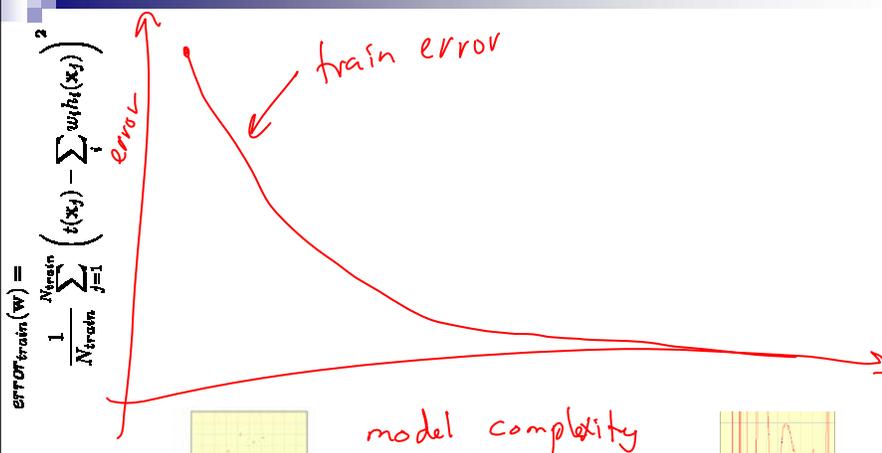
Training set error $\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_j \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$

- Given a dataset (Training data)
- Choose a loss function
 - e.g., squared error (L_2) for regression, *eg. 2, data likelihood*
- Training set error: For a particular set of parameters, loss function on training data:

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

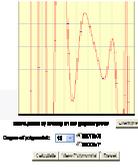
©Carlos Guestrin 2005-2007

Training set error as a function of model complexity



©Carlos Guestrin 2005-2007

Prediction error

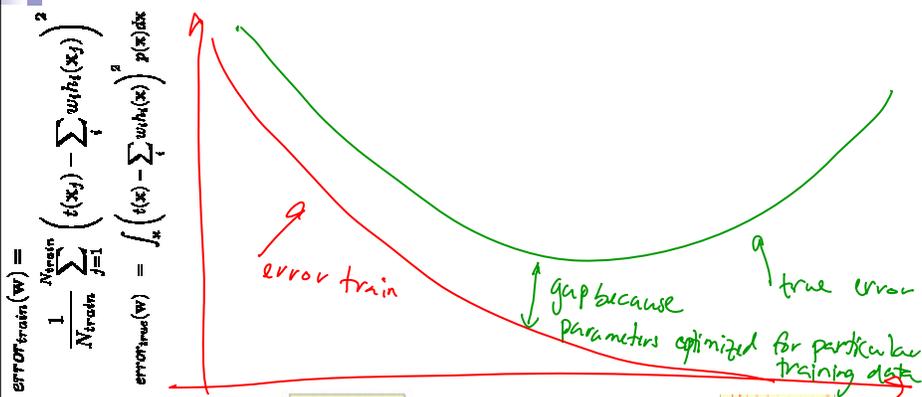


- Training set error can be poor measure of “quality” of solution
- Prediction error: We really care about error over all possible input points, not just training data:

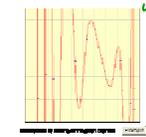
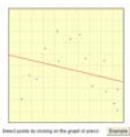
$$\begin{aligned}
 \text{error}_{\text{true}}(\mathbf{w}) &= E_{\mathbf{x}} \left[\left(t(\mathbf{x}) - \sum_i w_i h_i(\mathbf{x}) \right)^2 \right] \\
 &= \int_{\mathbf{x}} \left(t(\mathbf{x}) - \sum_i w_i h_i(\mathbf{x}) \right)^2 p(\mathbf{x}) d\mathbf{x}
 \end{aligned}$$

©Carlos Guestrin 2005-2007

Prediction error as a function of model complexity



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



©Carlos Guestrin 2005-2007

Computing prediction error

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

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

- 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

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

= error_{test}(w)

©Carlos Guestrin 2005-2007

Why training set error doesn't approximate prediction error?

- Sampling approximation of prediction error:

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

- Training error :

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

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

©Carlos Guestrin 2005-2007

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???

©Carlos Guestrin 2005-2007

Test set error

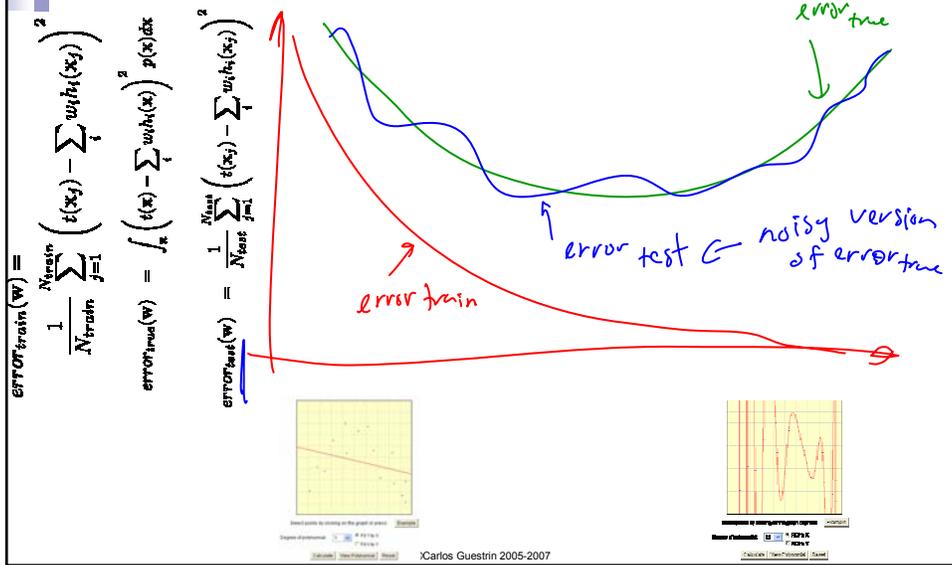
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_j \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

- 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:

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

©Carlos Guestrin 2005-2007

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:

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

- w was better in train than w' , but worse in true error.
- w' worse in train, but better in true than w .

How many points to I use for training/testing?

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

©Carlos Guestrin 2005-2007

Error estimators

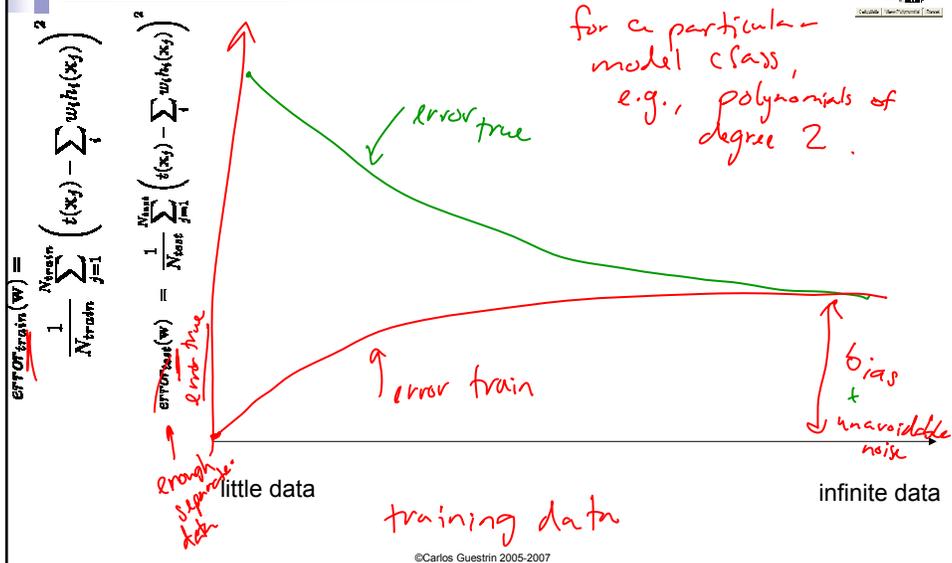
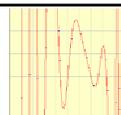
$$\text{error}_{\text{true}}(\mathbf{w}) = \int_{\mathbf{x}} \left(t(\mathbf{x}) - \sum_i w_i h_i(\mathbf{x}) \right)^2 p(\mathbf{x}) d\mathbf{x} \quad \leftarrow \text{gold standard}$$

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

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

©Carlos Guestrin 2005-2007

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 ~~ever~~ ~~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)

$$\text{error}_{\text{test}}(\mathbf{w}) = \frac{1}{N_{\text{test}}} \sum_{j=1}^{N_{\text{test}}} \left(t(x_j) - \sum_t w_t h_t(x_j) \right)^2$$

Announcements

- First homework is out:
 - Programming part and Analytic part
 - Remember collaboration policy: can discuss questions, but need to write your own solutions and code
 - Remember you are not allowed to look at previous years' solutions, search the web for solutions, use someone else's solutions, etc.
 - Due Oct. 3rd beginning of class
 - Start early!
- Recitation this week:
 - Bayes optimal classifiers, Naïve Bayes

©Carlos Guestrin 2005-2007

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, *data log likelihood for coin flips NBS*
- 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

©Carlos Guestrin 2005-2007

Types of (supervised) learning problems, revisited

- **Regression**, e.g., *sensor data modeling*
 - dataset: ⟨position; temperature⟩ *< 2,2 ; 23°C >*
 - hypothesis space: *poly degree 8* $H: \text{choose } h$
 - Loss function: *Squared loss* $h: X \mapsto \mathbb{R}$

- **Density estimation**, e.g., *curve for grades* $h: X \mapsto [0,1]$
 - dataset: ⟨grades⟩ *, 97, 95, 87, 99, ...* *s.t. $\sum_x h(x) = 1$*
 - hypothesis space: *Normal dist.* 
 - Loss function: *data likelihood* $h: X \mapsto \{1, 2, \dots, k\}$

- **Classification**, e.g., *Spam, reading your mind*
 - dataset: ⟨brain image; {verb v. noun}⟩
 - hypothesis space: *NB classifying & discrete (categorical) data*
 - Loss function: *data likelihood*

©Carlos Guestrin 2005-2007

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, *i.e., optimizes loss function*
- **Regression**
- **Density estimation**
- **Classification**
- (Not surprisingly) Seemly different problem, very similar solutions...

©Carlos Guestrin 2005-2007

What is NB really optimizing?

- Naïve Bayes assumption:

- Features are independent given class:

$$P(X_1, X_2 | Y) = P(X_1 | X_2, Y) P(X_2 | Y)$$

- More generally:

$$P(X_1 \dots X_n | Y) = \prod_i P(X_i | Y)$$

- NB Classifier:

$$P(Y, x_1, \dots, x_n) = P(Y) \prod_{i=1}^n P(x_i | Y)$$

©Carlos Guestrin 2005-2007

MLE for the parameters of NB

- Given dataset

- Count(A=a, B=b) ← number of examples where A=a and B=b

total data: m data point

- MLE for NB, simply:

- Prior: $P(Y=y) = \frac{\text{Count}(Y=y)}{m}$

- Likelihood: $P(X_i=x_i | Y=y) = \frac{\text{Count}(X_i=x_i, Y=y)}{\text{Count}(Y=y)}$

e.g., ith pixel is on given read a verb

the basin



©Carlos Guestrin 2005-2007

What is NB really optimizing?

Three independent binomial learning problems, each solved using counts

Let's use an example

MLE:

$$\max_{\theta} \ln P(D | \theta = \{\theta_{y_1}, \theta_{11t}, \theta_{11f}, \theta_{21t}, \theta_{21f}\}) =$$

loss data likelihood

$$\max_{\theta} \sum_{j=1}^m \ln P(y^{(j)}, x_1^{(j)}, x_2^{(j)} | \theta)$$

$$\max_{\theta} \sum_{j=1}^m \ln P(y^{(j)} | \theta_{y_1}) \cdot P(x_1^{(j)} | y^{(j)}, \theta_{11t}, \theta_{11f}) \cdot P(x_2^{(j)} | y^{(j)}, \theta_{21t}, \theta_{21f}) =$$

$$\sum_{j=1}^m \ln P(y^{(j)} | \theta_{y_1}) + \sum_{j=1}^m \ln P(x_1^{(j)} | y^{(j)}, \theta_{11t}, \theta_{11f}) + \sum_{j=1}^m \ln P(x_2^{(j)} | y^{(j)}, \theta_{21t}, \theta_{21f}) =$$

parameters to learn

$X_1 \leftarrow$ binary
 $X_2 \leftarrow$ binary
 $Y \leftarrow$ binary

$P(Y, X_1, X_2) = P(Y) \cdot P(X_1 | Y) \cdot P(X_2 | Y)$

$P(Y=t) = \theta_{y_1}$
 $P(Y=f) = 1 - \theta_{y_1}$

$P(X_1 | Y=t) = \theta_{11t}$
 $P(X_1 | Y=f) = \theta_{11f}$

$P(X_2 | Y=t) = \theta_{21t}$
 $P(X_2 | Y=f) = \theta_{21f}$

$\left[\max_{\theta_{y_1}} \sum_{j=1}^m \ln P(y^{(j)} | \theta_{y_1}) \right] + \left[\max_{\theta_{11t}, \theta_{11f}} \sum_{j=1}^m \ln P(x_1^{(j)} | y^{(j)}, \theta_{11t}, \theta_{11f}) \right] + \left[\max_{\theta_{21t}, \theta_{21f}} \sum_{j=1}^m \ln P(x_2^{(j)} | y^{(j)}, \theta_{21t}, \theta_{21f}) \right]$

©Carlos Guestrin 2005-2007

Generative v. Discriminative classifiers – Intuition

generate spm: sample (or set) $P(Y = \text{spm})$

sample words: $P(X | Y = \text{spm})$

- Want to Learn: $h: X \mapsto Y \leftarrow \{1, 2, 3, \dots, k\}$
 - X – features
 - Y – target classes
- Bayes optimal classifier – $P(Y|X)$
- Generative classifier, e.g., Naïve Bayes:
 - Assume some functional form for $P(X|Y), P(Y)$
 - Estimate parameters of $P(X|Y), P(Y)$ directly from training data
 - Use Bayes rule to calculate $P(Y|X=x) = \frac{P(Y, X=x)}{P(X=x)}$
 - This is a 'generative' model
 - Indirect computation of $P(Y|X)$ through Bayes rule
 - But, can generate a sample of the data, $P(X) = \sum_y P(y) P(X|y)$
- Discriminative classifiers, e.g., Logistic Regression:
 - Assume some functional form for $P(Y|X)$
 - Estimate parameters of $P(Y|X)$ directly from training data
 - This is the 'discriminative' model
 - Directly learn $P(Y|X)$
 - But cannot obtain a sample of the data, because $P(X)$ is not available

eg. NB: $P(X|Y) = \prod P(x_i | Y)$

of classification time: input x answer $P(Y|X=x)$

©Carlos Guestrin 2005-2007