

CS 2750 Machine Learning

Lecture 12a

Classification: Decision trees

Milos Hauskrecht

milos@cs.pitt.edu

5329 Sennott Square

Midterm exam

Midterm Wednesday, March 14, 2012

- **In-class (75 minutes)**
- **closed book**
- **material covered before Spring break**

Project proposals

Due: Monday, March 19, 2012

- **1 page long**

Proposal

- **Written proposal:**
 1. Outline of a learning problem, type of data you have available. Why is the problem important?
 2. Learning methods you plan to try and implement for the problem. References to previous work.
 3. How do you plan to test, compare learning approaches
 4. Schedule of work (approximate timeline of work)

Project proposals

Where to find the data:

- From your research
- UC Irvine data repository
- Various text document repositories
- I have some bioinformatics data I can share but other data can be found on the NIH or various university web sites (e.g. microarray data, proteomic data)
- Synthetic data that are generated to demonstrate your algorithm works

Project proposals

Problems to address:

- Get the ideas for the project by browsing the web
- It is tempting to go with simple classification but definitely try to add some complexity to your investigations
- Multiple, not just one method, try some more advanced methods, say those that combine multiple classifiers to learn a model (ensemble methods) or try to modify the existing methods

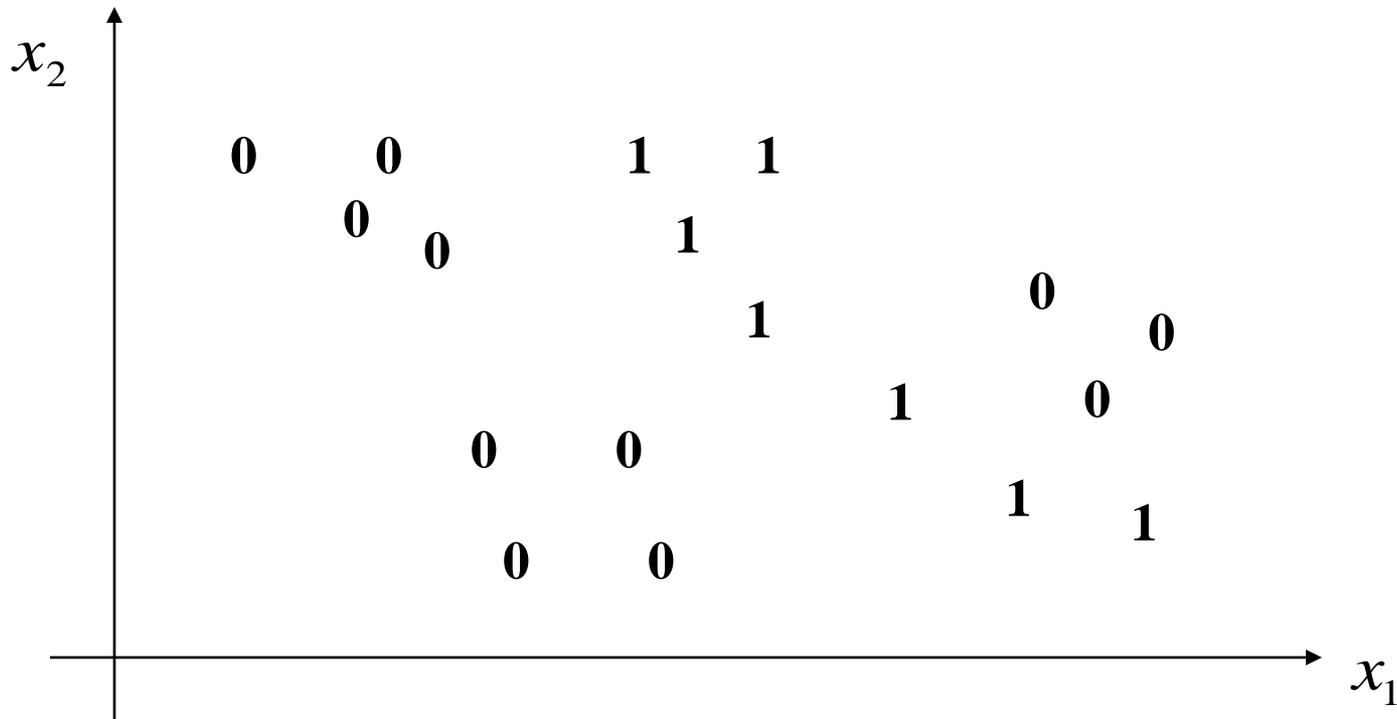
Project proposals

Interesting problems to consider:

- Advanced methods for learning multi-class problems
- Learning the parameters and structure of Bayesian Belief networks
- Clustering of data – how to group examples
- Dimensionality reduction/feature selection – how to deal with a large number of inputs
- Learning how to act – Reinforcement learning
- Anomaly detection – how to identify outliers in data

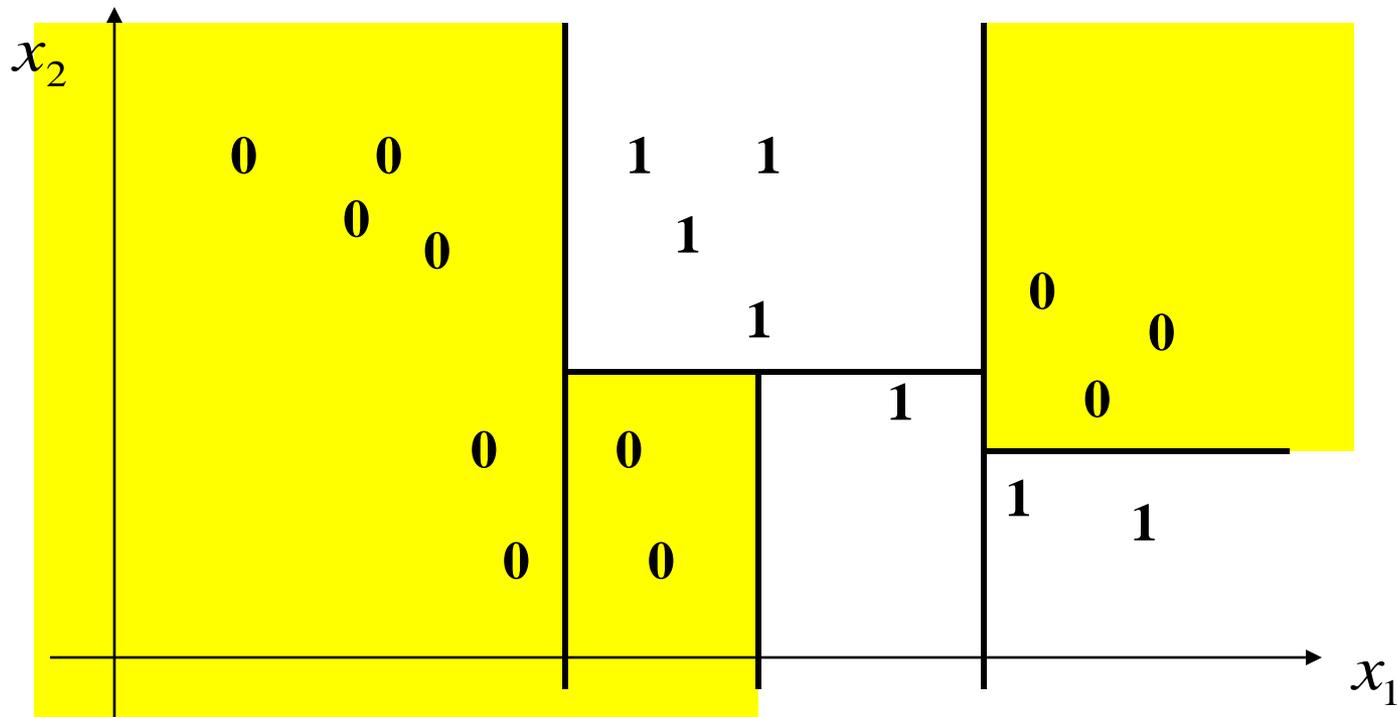
Decision trees

- An alternative approach to classification:
 - **Partition the input space to regions**
 - **Regress or classify independently in every region**



Decision trees

- An alternative approach to classification:
 - **Partition the input space to regions**
 - **Regress or classify independently in every region**



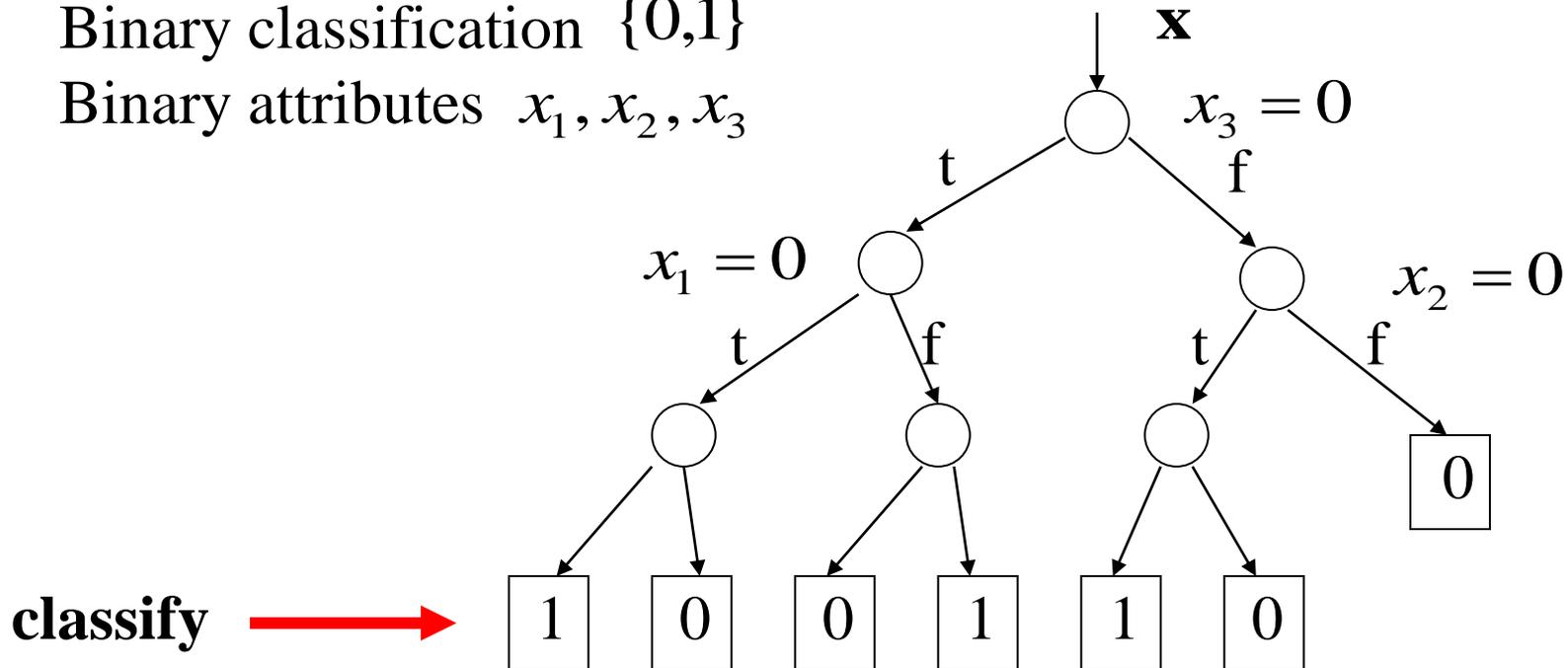
Decision trees

- The partitioning idea is used in the **decision tree model**:
 - Split the space recursively according to inputs in \mathbf{x}
 - Regress or classify at the bottom of the tree

Example:

Binary classification $\{0,1\}$

Binary attributes x_1, x_2, x_3



Decision trees

How to construct the decision tree?

- **Top-bottom algorithm:**

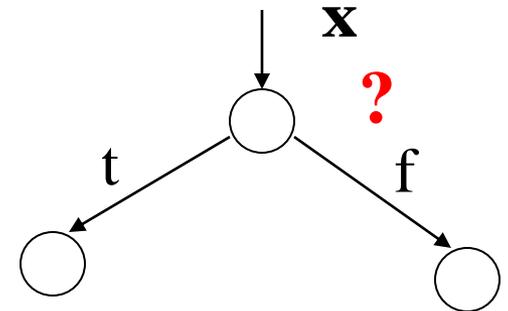
- Find the best split condition (quantified based on the impurity measure)
- Stops when no improvement possible

- **Impurity measure:**

- Measures how well are the two classes separated
- Ideally we would like to separate all 0s and 1

- Splits of **finite vs. continuous value attributes**

Continuous value attributes conditions: $x_3 \leq 0.5$



Impurity measure

Let $|D|$ - Total number of data entries

$|D_i|$ - Number of data entries classified as i

$p_i = \frac{|D_i|}{|D|}$ - ratio of instances classified as i

- **Impurity measure** defines how well the classes are separated
- In general the impurity measure should satisfy:
 - Largest when data are split evenly for attribute values

$$p_i = \frac{1}{\text{number of classes}}$$

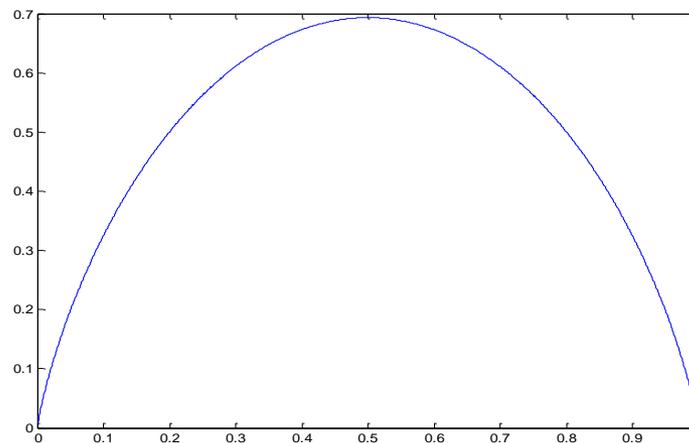
- Should be 0 when all data belong to the same class

Impurity measures

- There are various impurity measures used in the literature
 - **Entropy based measure (Quinlan, C4.5)**

$$I(D) = \text{Entropy}(D) = -\sum_{i=1}^k p_i \log p_i$$

Example for k=2



- **Gini measure (Breiman, CART)**

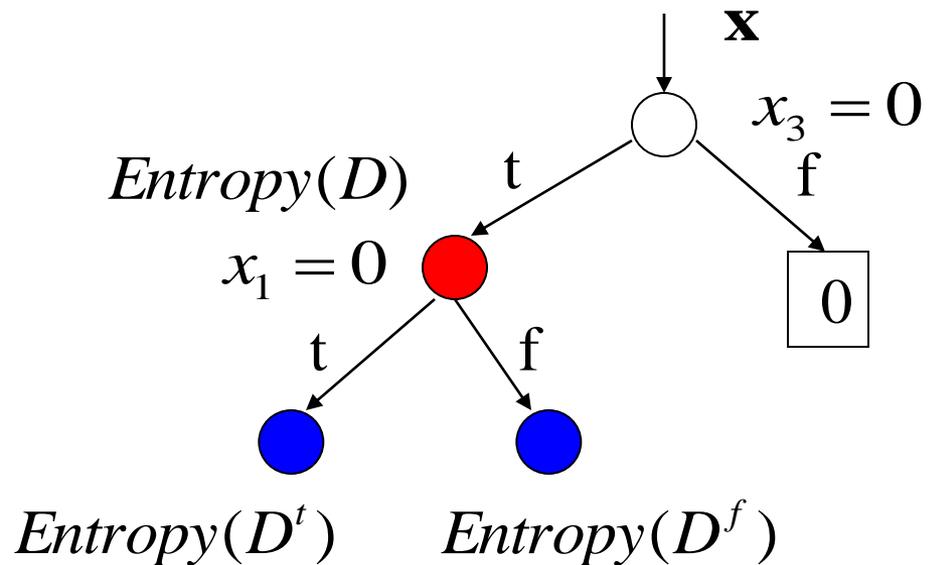
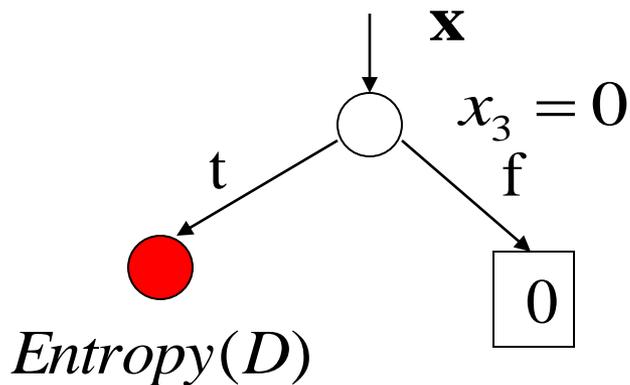
$$I(D) = \text{Gini}(D) = 1 - \sum_{i=1}^k p_i^2$$

Impurity measures

- **Gain due to split** – expected reduction in the impurity measure (entropy example)

$$Gain(D, A) = Entropy(D) - \sum_{v \in Values(A)} \frac{|D^v|}{|D|} Entropy(D^v)$$

$|D^v|$ - a partition of D with the value of attribute $A = v$



Decision tree learning

- **Greedy learning algorithm:**

- Repeat until no or small improvement in the purity

- Find the attribute with the highest gain

- Add the attribute to the tree and split the set accordingly

- Builds the tree in the top-down fashion

- Gradually expands the leaves of the partially built tree

- The method is greedy

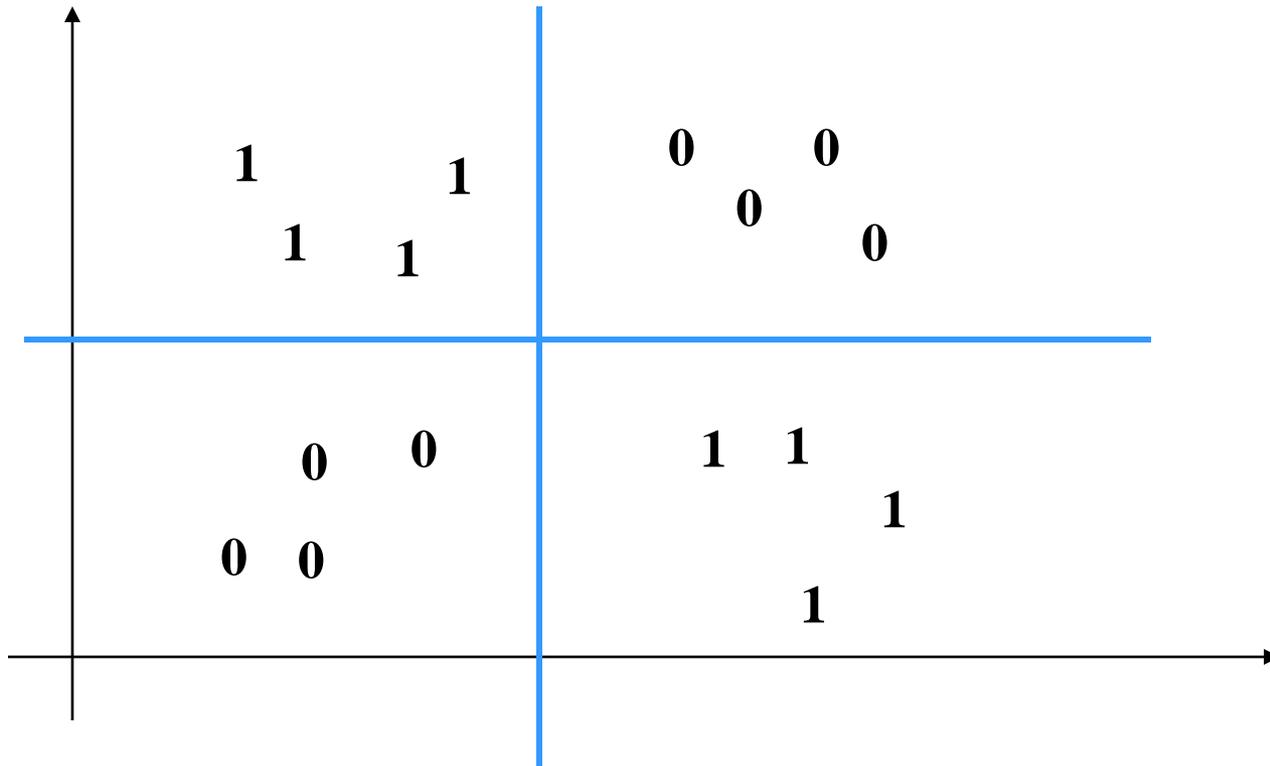
- It looks at a single attribute and gain in each step

- May fail when the combination of attributes is needed to improve the purity (parity functions)

Decision tree learning

- **Limitations of greedy methods**

Cases in which a combination of two or more attributes improves the impurity



Decision tree learning

By reducing the impurity measure we can grow **very large trees**

Problem: Overfitting

- We may split and classify very well the training set, but we may do worse in terms of the generalization error

Solutions to the overfitting problem:

- **Solution 1.**
 - Prune branches of the tree built in the first phase
 - Use validation set to test for the overfit
- **Solution 2.**
 - Test for the overfit in the tree building phase
 - Stop building the tree when performance on the validation set deteriorates