# CSDS 440: Machine Learning

## Soumya Ray (he/him, sray@case.edu)

Olin 516

Office hours T, Th 11:15-11:45 or by appointment

# Recap

- A learning system is specified by a g_____, task e_____ and a p_____ m_____.
- What are the two phases of learning? What happens in these phases?
- What are online and offline learning?
- Every ML system must reason from s_____ to the g_____ c____. This is called i_____ g_____.
- The system is looking for the t_____ c_____, which is the _____.
- To find this the system searches a h_____ s_____.
- All possible hypotheses can/cannot be considered. Why?
- This is called _____.
- What is the "inductive bias" of a learning algorithm?

# Today

- Foundations of machine Learning

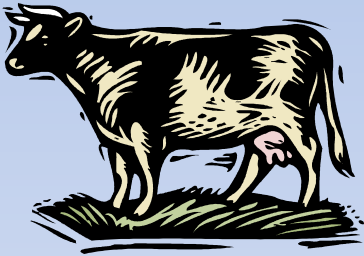Soumya Ray, Case Western Reserve U.

# Supervised Learning

- Examples *E* are annotated with target concept's output by a teacher/oracle

- Learning system must find a concept that matches annotations (*P*)
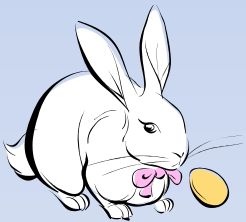
- Example: learn to recognize animals

# Supervised Learning



tiger

cow

elephant

Note: Annotation received by learner does not need to be correct!!

starfish

# Other Learning Paradigms

- Unsupervised Learning
- Semi-supervised Learning
- Active Learning
- Transductive Learning
- Transfer Learning
- Structured Prediction
- Reinforcement Learning
- Preference Learning (Ranking)
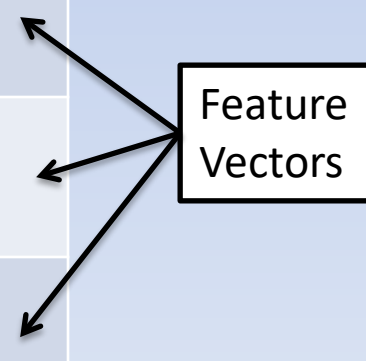- "Few-shot" learning

# Example Representation

- What is the *internal representation* of an example in a learning system?

- Representation choice affects reasoning and the choice of hypothesis space, and the cost of learning

# Feature Vector Representation

- Examples are attribute-value pairs (note "feature"=="attribute")
- Number of attributes are fixed
- Can be written as an $n$-by-$m$ matrix

|  | Attribute$_1$ | Attribute$_2$ | Attribute$_3$ |
|---|---|---|---|
| **Example$_1$** | Value$_{11}$ | Value$_{12}$ | Value$_{13}$ |
| **Example$_2$** | Value$_{21}$ | Value$_{22}$ | Value$_{23}$ |
| **Example$_3$** | Value$_{31}$ | Value$_{32}$ | Value$_{33}$ |

Feature Vectors

# Example

| | Has-fur? | Long-Teeth? | Scary? |
|---|---|---|---|
| **Animal$_1$** | Yes | No | No |
| **Animal$_2$** | No | Yes | Yes |
| **Animal$_3$** | Yes | Yes | Yes |

# Types of Features

- Discrete, Nominal

- Continuous

- Discrete, Ordered

- Hierarchical

- *Color $\in$ (red, blue, green)*

- *Height*

- *Size $\in$ (small, medium, large)*

- *Shape $\in$* **closed**

  **polygon**     **continuous**

  **square**   **triangle**   **circle**    **ellipse**

# Feature Space

- We can think of examples embedded in an $n$ dimensional vector space

# Other Example Representations

- Relational representation

- Multiple-instance representation

- Sequential representation
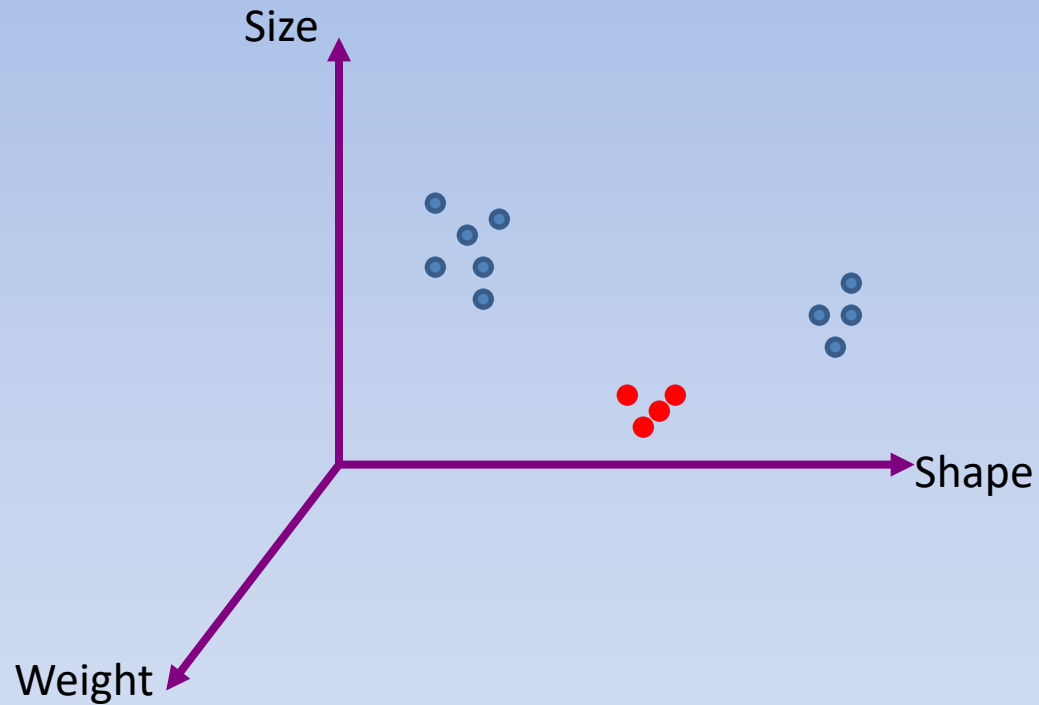
- Multi-view representation

# The Binary Classification Problem

- Simplest supervised learning problem

- Target concept assigns one of two labels ("*positive*" or "*negative*") to all examples---the <span style="color:red">class label</span>

- Can extend to "multiclass", "regression", "multi-label" problems

# Example

| | Has-fur? | Long-Teeth? | Scary? | *Lion?* |
|---|---|---|---|---|
| **Animal$_1$** | Yes | No ($x_{ij}$) | No | No |
| **Animal$_2$** | No | Yes | Yes | No |
| **Animal$_3$** | Yes | Yes | Yes | Yes |

$X$ spans Has-fur?, Long-Teeth?, Scary?. $Y$ spans Lion?
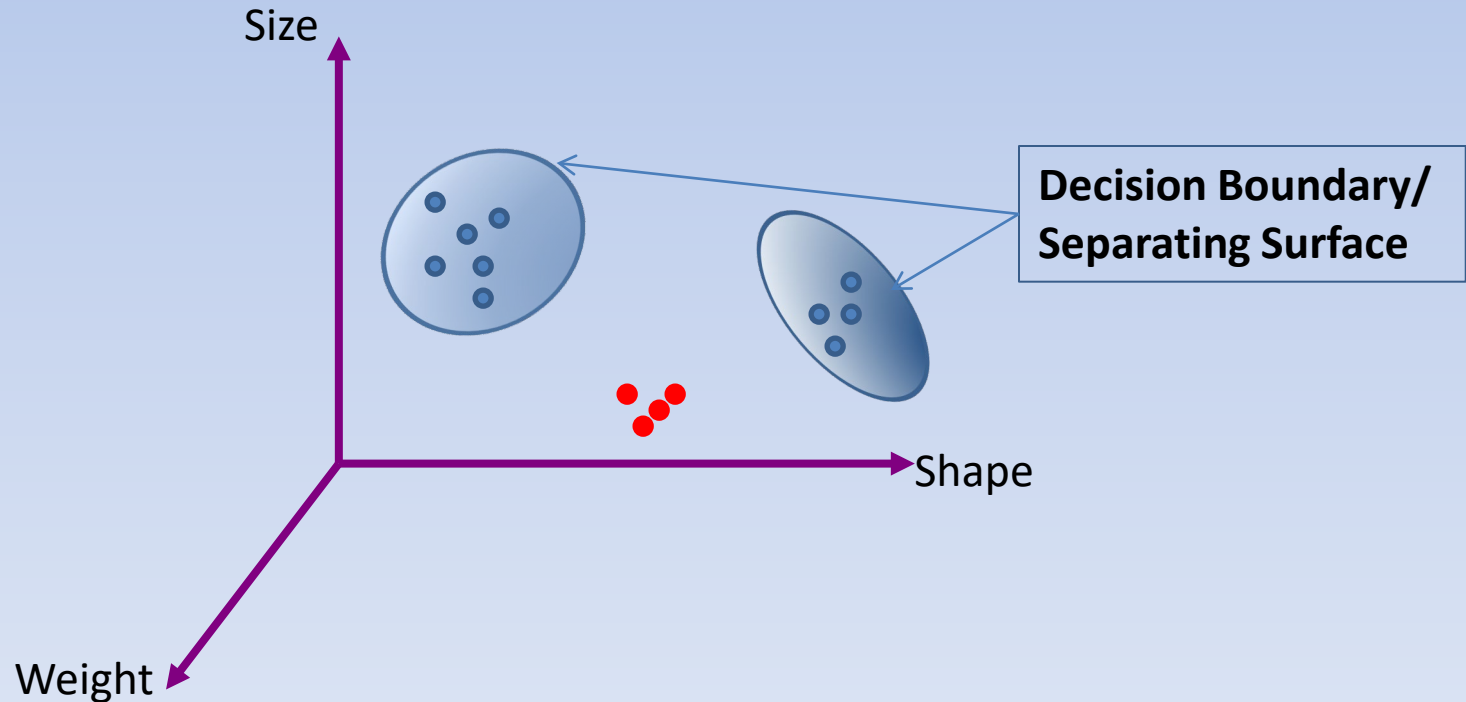
$(\boldsymbol{x}_i,\ y_i)$

# Example in Feature Space

# The Learning Problem

- Given: A binary classification problem

- Do: Produce a "classifier" (concept) that assigns a label to a new example

# Binary Classifier Concept Geometry

- (Union of ) $N$-dimensional volume(s) in feature space (possibly a disjoint collection)



Decision Boundary/ Separating Surface

# Decision Tree Induction (Ch 3, Mitchell)

- A "classical" (1980s) family of machine learning algorithms for classification


- Widely used and extremely popular, available in nearly all ML toolkits
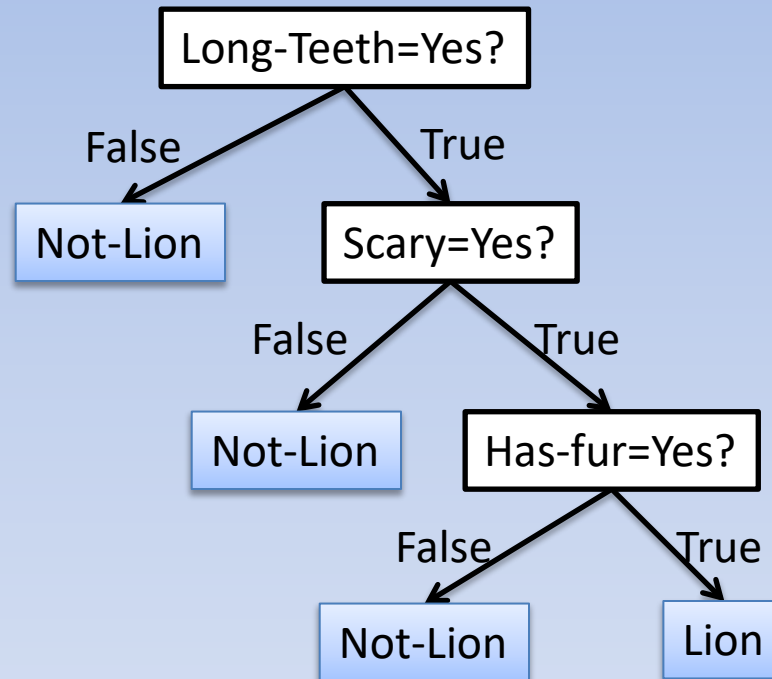
# What is a Decision Tree?

- Tree: directed acyclic graph, each node has at most one parent

- Internal nodes: Tests on attributes

- Leaves: Class labels

# Example

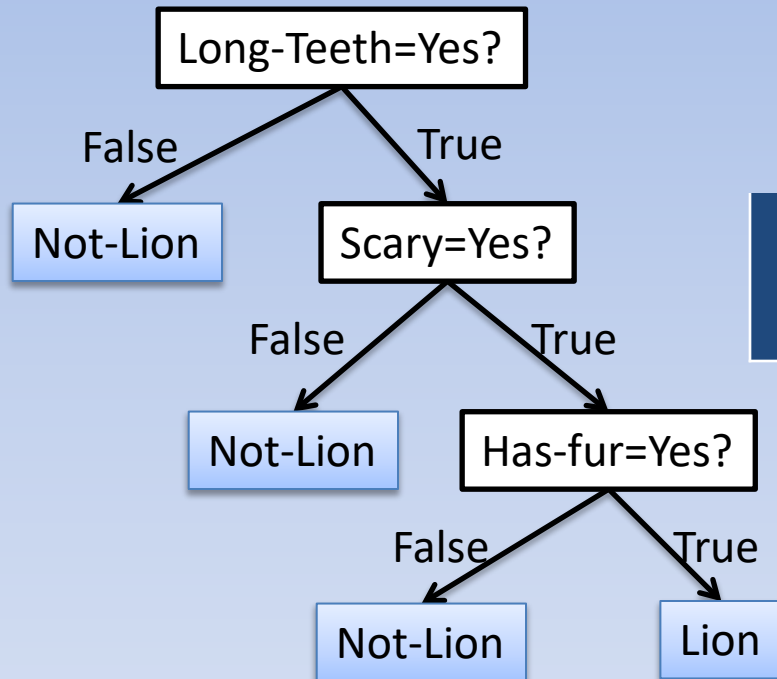| | Has-fur? | Long-Teeth? | Scary? | *Lion?* |
|---|---|---|---|---|
| **Animal$_1$** | Yes | No | No | No |
| **Animal$_2$** | No | Yes | Yes | No |
| **Animal$_3$** | Yes | Yes | Yes | Yes |

# Example

# Classification with a decision tree

- Suppose we are given a tree and a new example

- Starting at the root, check each attribute test

- This identifies a path through the tree, follow this until we reach a leaf

- Assign the class label in the leaf

# Example



Long-Teeth=Yes?
- False → Not-Lion
- True → Scary=Yes?
  - False → Not-Lion
  - True → Has-fur=Yes?
    - False → Not-Lion
    - True → Lion

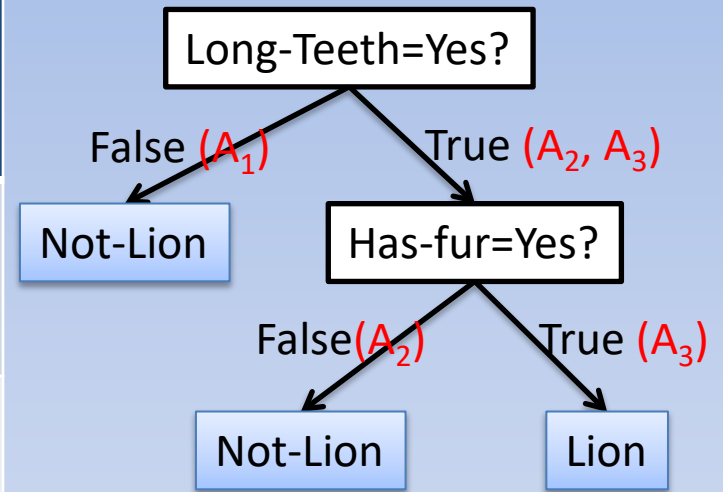| | Has-fur? | Long-Teeth? | Scary? |
|---|---|---|---|
| **Animal$_1$** | Yes | Yes | No |

# Decision Tree Induction

- Given a set of examples, produce a decision tree
- Decision tree induction works using the idea of <span style="color:red">recursive partitioning</span>
  - At each step, the algorithm will <span style="color:red">choose an attribute test</span>
    - If no attribute looks good, return
  - The chosen test will partition the examples into disjoint partitions
  - The algorithm will then recursively call itself on each partition until
    - a partition only has data from one class (<span style="color:red">pure</span> node) OR
    - it runs out of attributes

# Example

| | Has-fur? | Long-Teeth? | Scary? | *Lion?* |
|---|---|---|---|---|
| **Animal$_1$** | Yes | No | No | No |
| **Animal$_2$** | No | Yes | Yes | No |
| **Animal$_3$** | Yes | Yes | Yes | Yes |

Long-Teeth=Yes?

False ($A_1$)     True ($A_2$, $A_3$)

Not-Lion     Has-fur=Yes?

False($A_2$)     True ($A_3$)

Not-Lion     Lion

# Choosing an Attribute

- Which attribute should we choose to test first?

  – Ideally, the one that is "most predictive" of the class label

    - i.e., the one that gives us the "most information" about what the label should be

- This idea is captured by the "(Shannon) entropy" of a random variable

# Entropy of a Random Variable

- Suppose a random variable $X$ has density $p(x)$. Its (Shannon) "entropy" is defined by:

$$H(X) = E(-\log_2(p(X)))$$

$$= -\sum_x p(X = x) \log_2(p(X = x))$$

- Note: 0log(0) = 0 .