# CSDS 440: Machine Learning

Soumya Ray (he/him, [sray@case.edu](mailto:sray@case.edu))

Olin 516

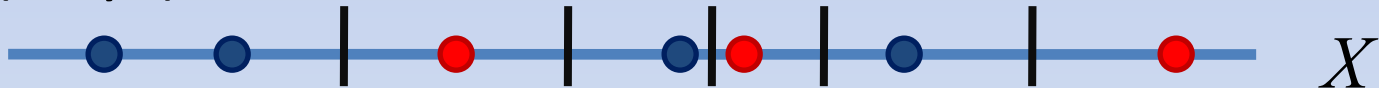Office hours T, Th 11:15-11:45 or by appointment

# Recap

- To choose a test, we look for an attribute that provides i_____ about the l____. A quantity that encodes this is the e_____ of a random variable.

- E_____ is the expected length of the s____ l____ d____ of a random variable.

- I____ g____ is the r____ of e____ of the class variable b___ and a____ partitioning.

- What problem arises with nominal features and info gain?

- We can attempt to resolve this issue by adjusting the split criterion. GR(X)=____/____. This works because _____.

- We partition on continuous features by considering all tests of the form _____.

- We only need to consider values that _____.

# Today

- Decision Tree Induction (Ch 3, Mitchell)
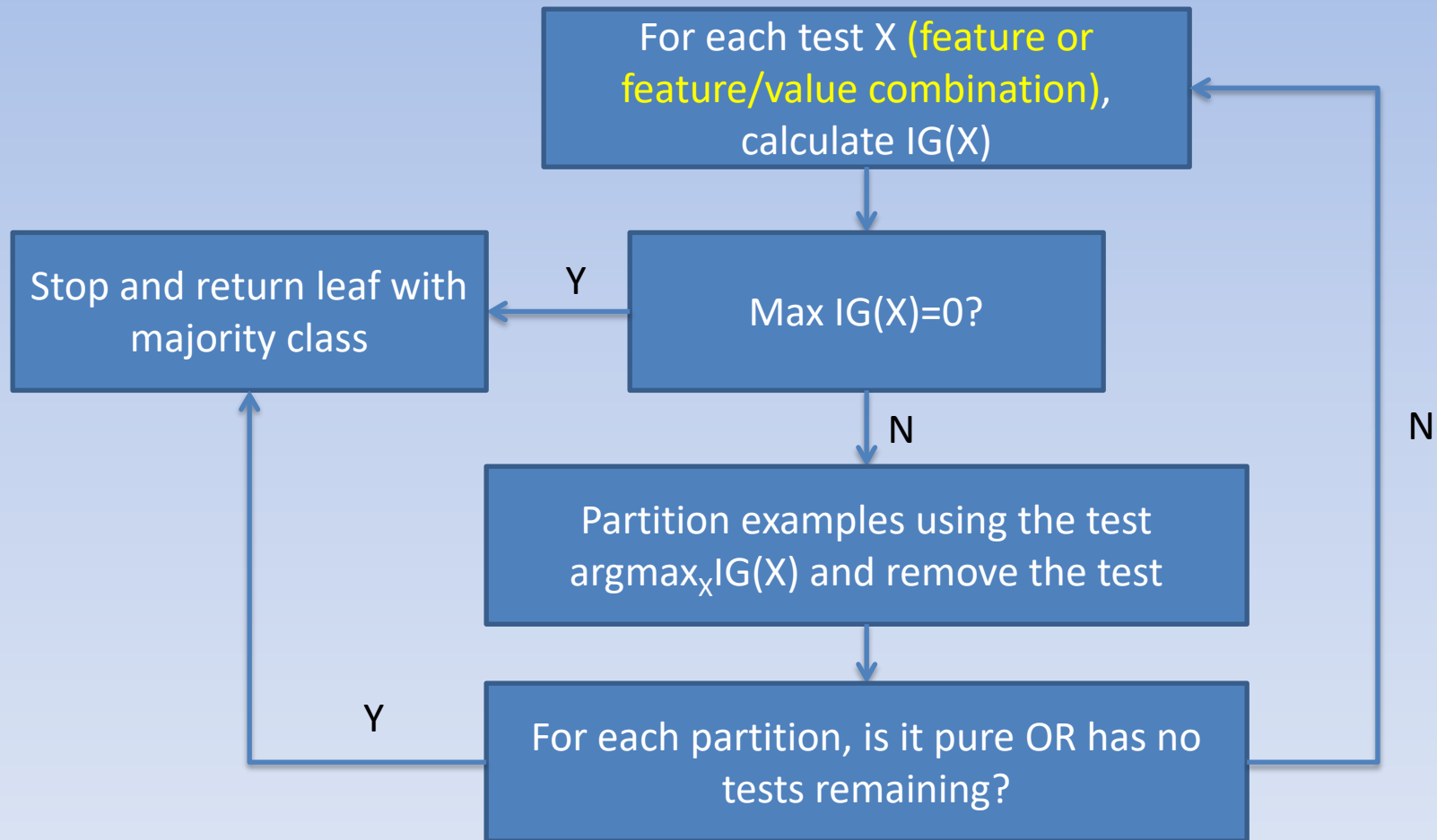- Overfitting and overfitting control

# Continuous Attributes

- Cannot test for equality
- Consider all Boolean tests of the form $X \geq v$ (or $X \leq v$)
  - Only values of interest are those $v$ that separate adjacent training examples with different classes (why?)

$$X$$

- Note: In this case, the attribute cannot be removed, though the test ((attribute, value) tuple) can be

# ID3 Algorithm---Training phase

For each test X (feature or feature/value combination), calculate IG(X)

Max IG(X)=0?

Y → Stop and return leaf with majority class

N

Partition examples using the test $\text{argmax}_X \text{IG}(X)$ and remove the test

For each partition, is it pure OR has no tests remaining?

Y

N

# Example

| Color | Area | Shape | Class Label |
|-------|------|-------|-------------|
| red | 0.1 | circle | 1 |
| blue | 0.2 | triangle | 1 |
| green | 0.3 | triangle | 0 |
| green | 0.3 | circle | 0 |
| green | 0.4 | square | 0 |
| red | 0.4 | triangle | 1 |
| blue | 0.6 | circle | 0 |
| red | 0.7 | square | 0 |
| blue | 0.8 | square | 0 |

# Example

| Color | Area | Shape | Class Label |
|-------|------|-------|-------------|
| ~~red~~ | ~~0.1~~ | ~~circle~~ | ~~1~~ |
| ~~blue~~ | ~~0.2~~ | ~~triangle~~ | ~~1~~ |
| green | 0.3 | triangle | 0 |
| green | 0.3 | circle | 0 |
| green | 0.4 | square | 0 |
| red | 0.4 | triangle | 1 |
| blue | 0.6 | circle | 0 |
| red | 0.7 | square | 0 |
| blue | 0.8 | square | 0 |

# Example

| Color | Area | Shape | Class Label |
|-------|------|-------|-------------|
| ~~red~~ | ~~0.1~~ | ~~circle~~ | ~~1~~ |
| ~~blue~~ | ~~0.2~~ | ~~triangle~~ | ~~1~~ |
| ~~green~~ | ~~0.3~~ | ~~triangle~~ | ~~0~~ |
| ~~green~~ | ~~0.3~~ | ~~circle~~ | ~~0~~ |
| ~~green~~ | ~~0.4~~ | ~~square~~ | ~~0~~ |
| red | 0.4 | triangle | 1 |
| ~~blue~~ | ~~0.6~~ | ~~circle~~ | ~~0~~ |
| red | 0.7 | square | 0 |
| ~~blue~~ | ~~0.8~~ | ~~square~~ | ~~0~~ |

# An Issue

- Given enough features, ID3 will usually be able to fit training examples exactly (i.e. every leaf is pure), because the tree can be grown as much as needed
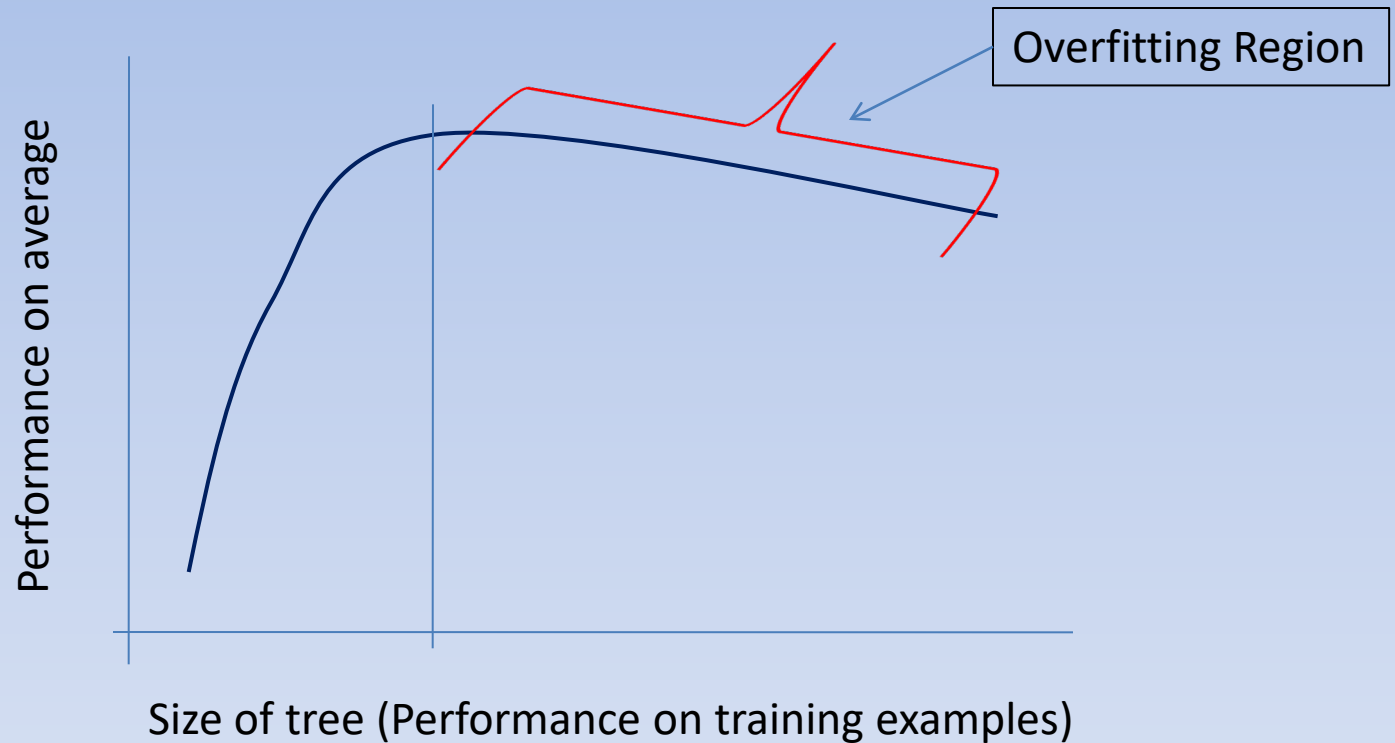
- But real data is noisy

# Overfitting



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

| | Has-fur? | Long-Teeth? | Scary? | Color? | *Lion?* |
|---|---|---|---|---|---|
| Animal₁ | Yes | No | No | Green | No |
| Animal₂ | No | Yes | No | Black | No |
| Animal₃ | Yes | Yes | Yes | Golden | Yes |
| Animal₄ | Yes | Yes | No | Blue | Yes |
| Animal₅ | Yes | Yes | Yes | Tawny | Yes |

# Overfitting

- If a learned concept $h$ has
  - Higher performance (lower error) on the training examples, BUT
  - Lower performance (higher error) on average across all examples

than some alternative concept $h'$ in the same hypothesis space, $h$ is said to have overfit to the training examples

# Overfitting



Overfitting Region

Performance on average

Size of tree (Performance on training examples)

# Controlling Overfitting

- Introduce a <span style="color:red">restriction</span> on the hypothesis space to prevent overly-complex hypotheses from being learned

  - Early Stopping
  - Post Pruning

# Early Stopping

- Standard algorithm stops growing the tree when $IG(X)=0$ for all $X$

- Early stopping stops growing the tree when $IG(X) \leq \varepsilon,$ for some chosen $\varepsilon$

- Sensitive to choice of $\varepsilon$

- Easy to implement, but does not work very well in practice
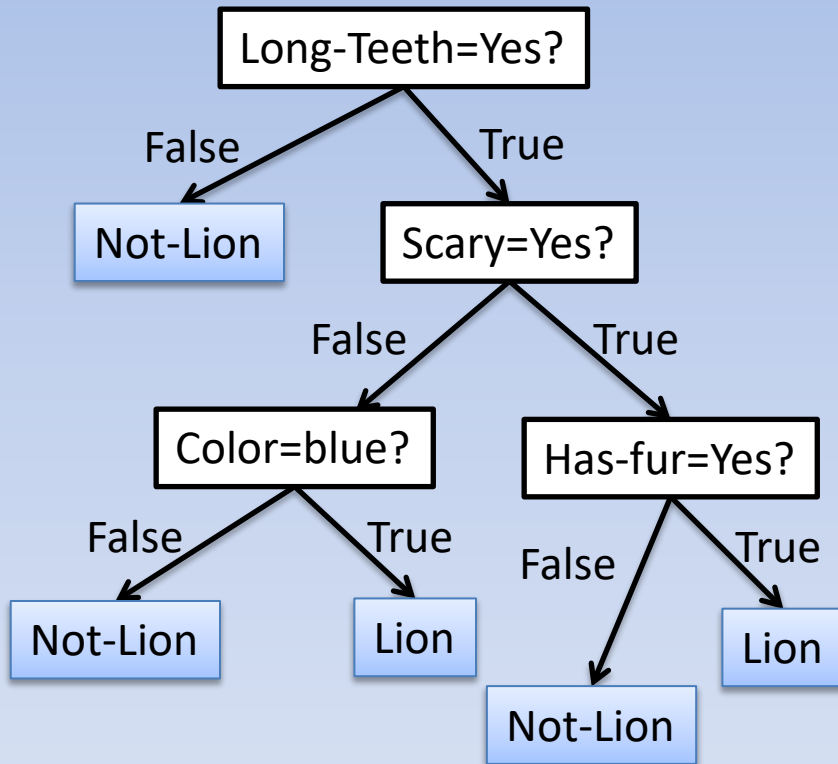
# Greedy post-pruning

- Hold aside some training examples at start (<span style="color:red">validation set</span>)

- Grow tree as usual on remainder

- Then run a *greedy pruning* algorithm

# Greedy post-pruning
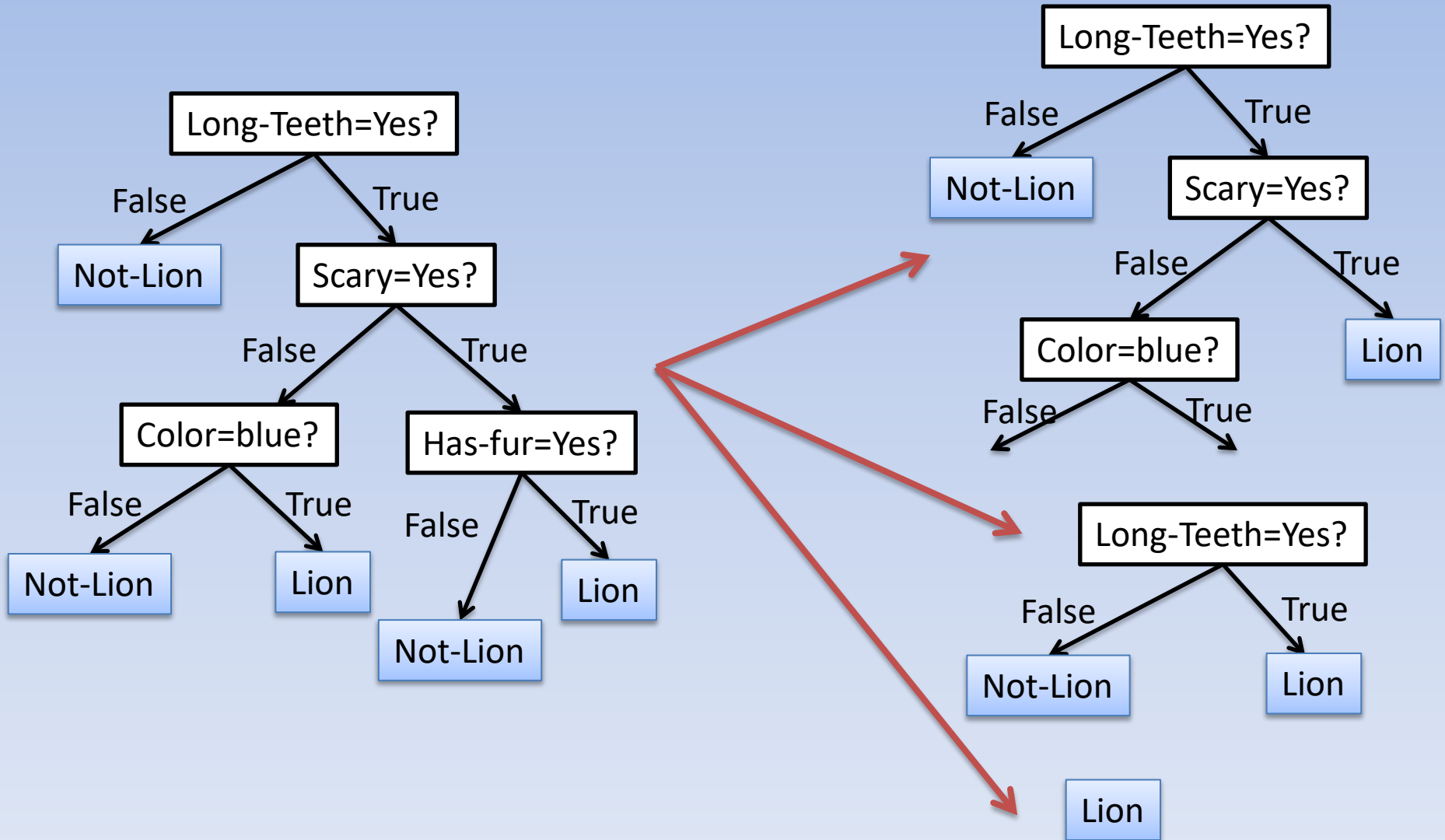
- For each internal node, construct a tree without that node

  – Convert node to leaf by predicting majority class

  – Delete subtree below node

- Evaluate this *pruned* tree <span style="color:red">on the validation set</span>

- Find the single node that improves performance the most over the unpruned tree and remove it

- Repeat steps above until no node removal improves performance

# Greedy Post Pruning

Long-Teeth=Yes?

False → Not-Lion

True → Scary=Yes?

Scary=Yes?
- False → Color=blue?
- True → Has-fur=Yes?

Color=blue?
- False → Not-Lion
- True → Lion

Has-fur=Yes?
- False → Not-Lion
- True → Lion

| | Has-fur? | Long-Teeth? | Scary? | Color? | *Lion?* |
|---|---|---|---|---|---|
| **Animal$_1$** | Yes | No | No | Green | No |
| **Animal$_2$** | No | Yes | No | Black | No |
| **Animal$_3$** | Yes | Yes | Yes | Golden | Yes |
| **Animal$_4$** | Yes | Yes | No | Blue | Yes |
| **Animal$_5$** | Yes | Yes | Yes | Tawny | Yes |
| **Animal$_6$** | No | Yes | No | Blue | No |

# Greedy Post Pruning

# Greedy Post Pruning



**Left tree:**

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

**Right tree:**

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