

# Decision Trees - Pruning

CS1090A Introduction to Data Science

Pavlos Protopapas, Kevin Rader, and Chris Gumb



Photo: Eleonore Wen  
Jungfrau Summit, Switzerland

**SO MY DAUGHTER SAYS  
YOU WORK IN DATA SCIENCE**

**YES, TODAY I BUILT  
A MODEL WITH  $R^2$  OF 1**

**YOU HAVE EXACTLY 10  
SECONDS TO GET OUT OF MY HOUSE**

# Outline

---

- Decision Trees – Regression
- Numerical vs Categorical Attributes
- Pruning

# Alternatives to Using Stopping Conditions

What is the major issue with pre-specifying a stopping condition?

- You may stop too early or stop too late.

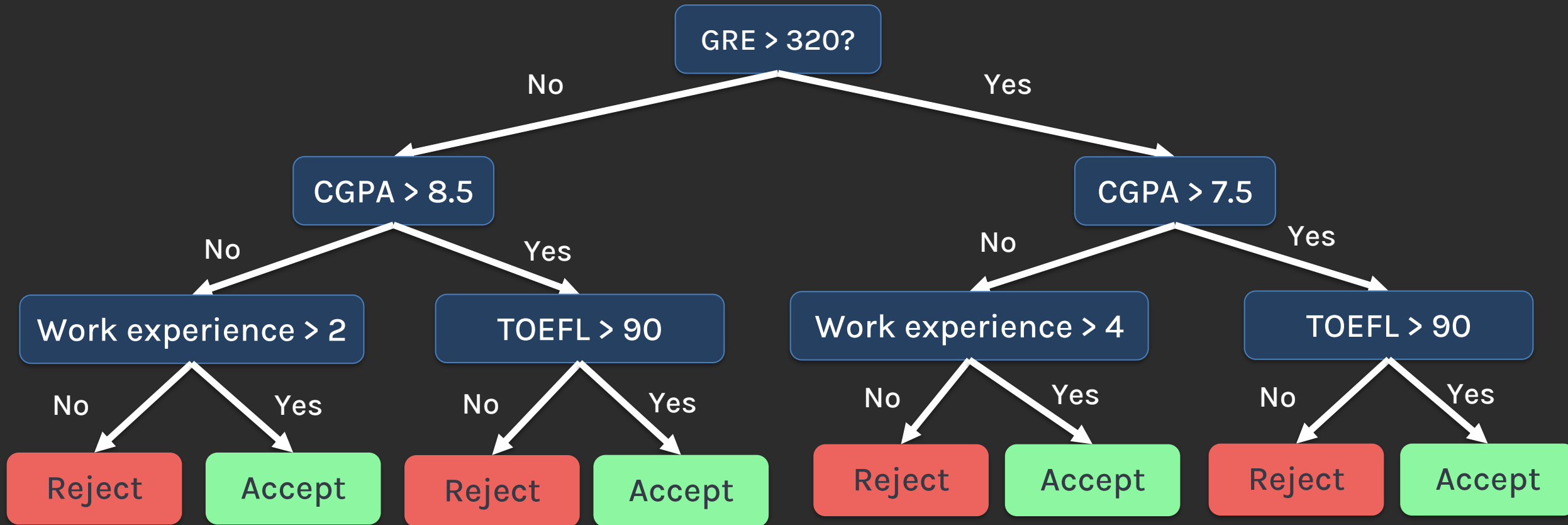
How can we fix this issue?

- Choose several stopping criteria (e.g., set minimal  $\text{Gain}(R)$  at various levels) and cross-validate to decide which one is the best.

What is an alternative approach to this method?

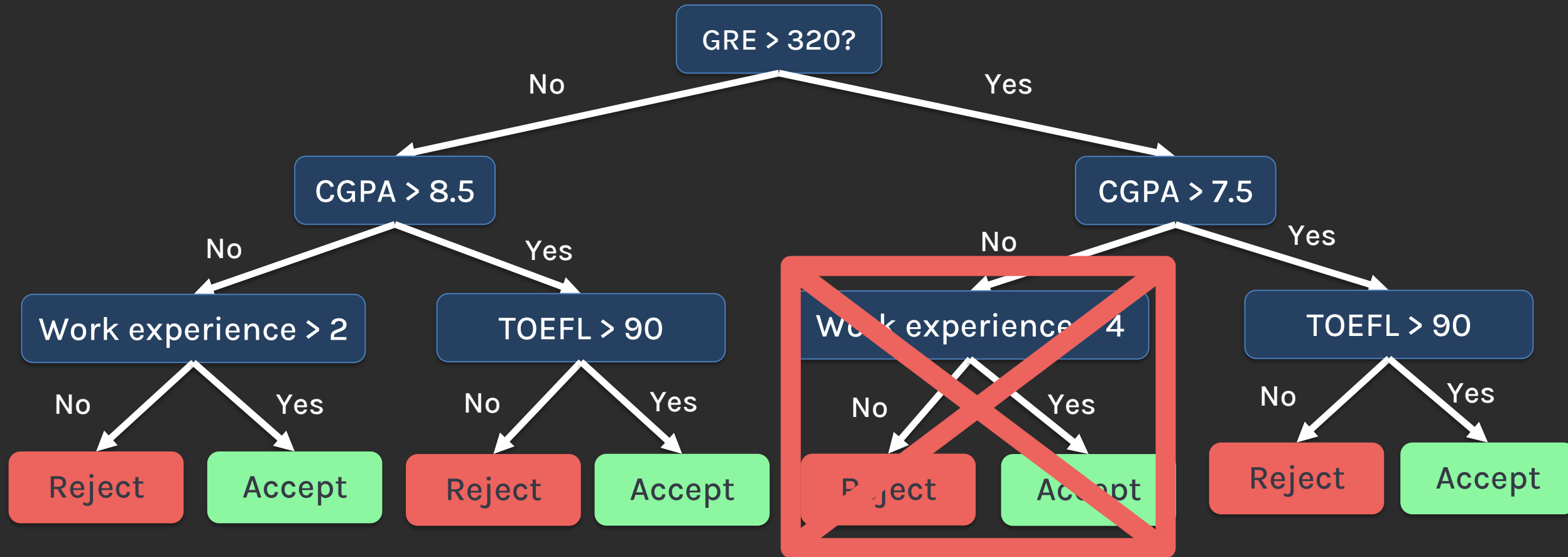
- Don't stop. Instead, prune the tree!

# Example: Evaluating Applications to a GRADUATE Program



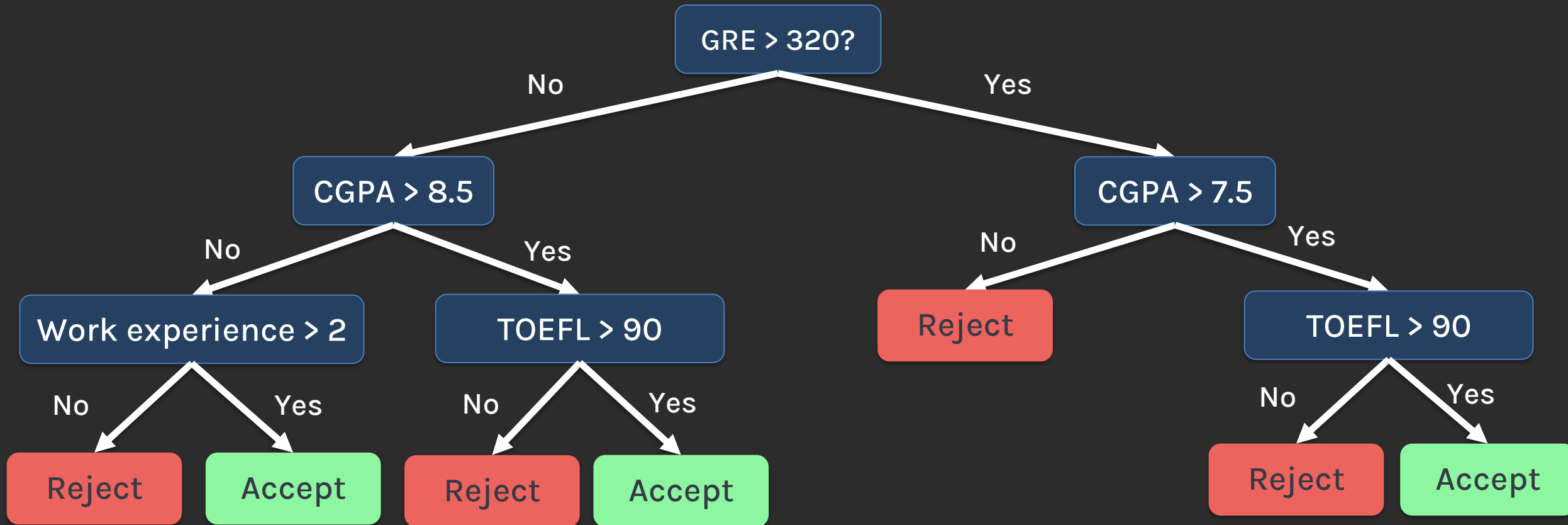
\*\* Above representation is only for pedagogical purposes.

# Example: Evaluating Applications to a GRADUATE Program



\*\* Above representation is only for pedagogical purposes.

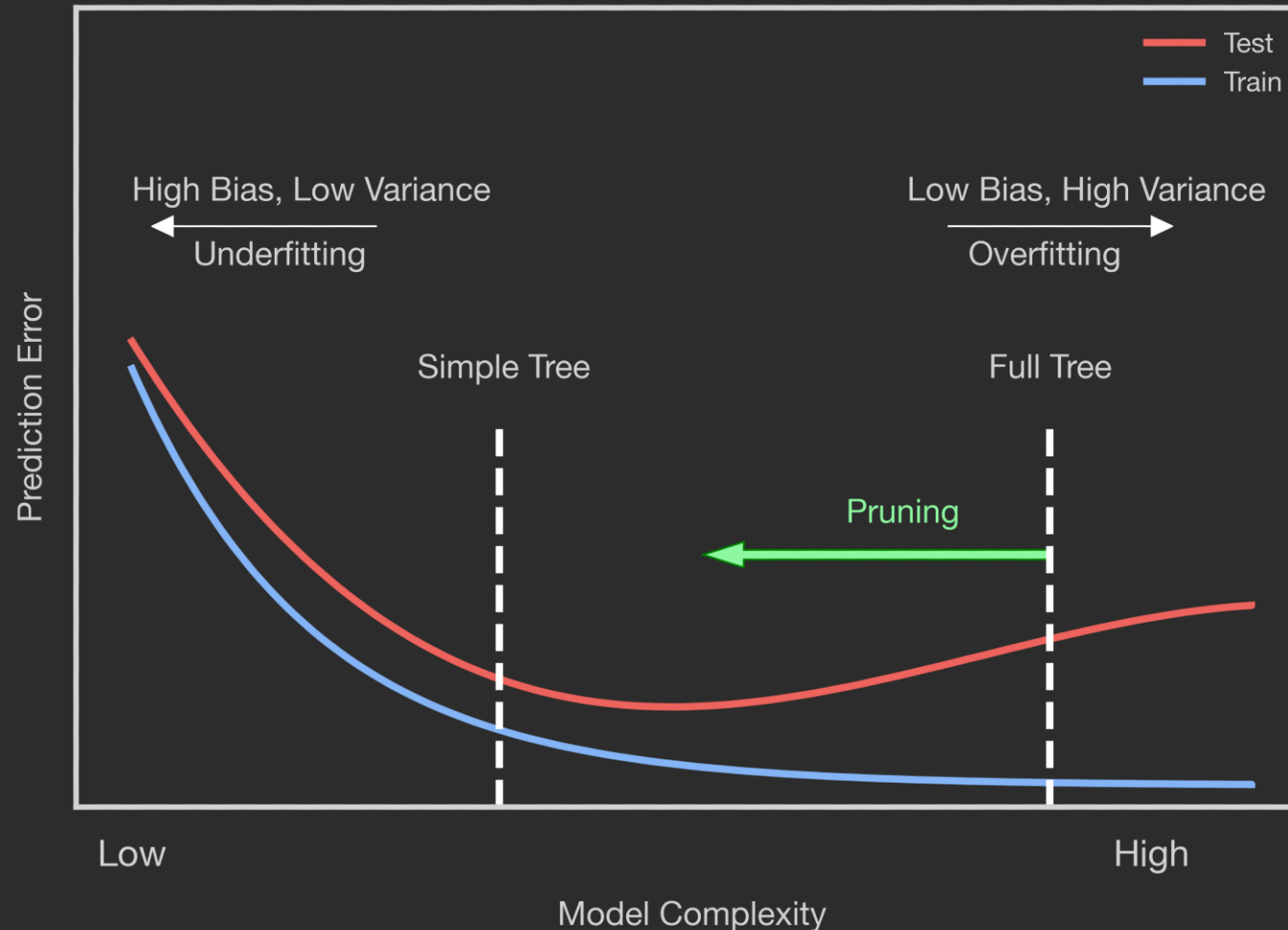
# Example: Evaluating Applications to a GRADUATE Program



\*\* Above representation is only for pedagogical purposes.



# Motivation for Pruning



Rather than preventing a complex tree from growing, we can obtain a simpler tree by ‘pruning’ a complex one.



# Pruning

There are many methods of pruning. A common one is the **cost complexity pruning**:

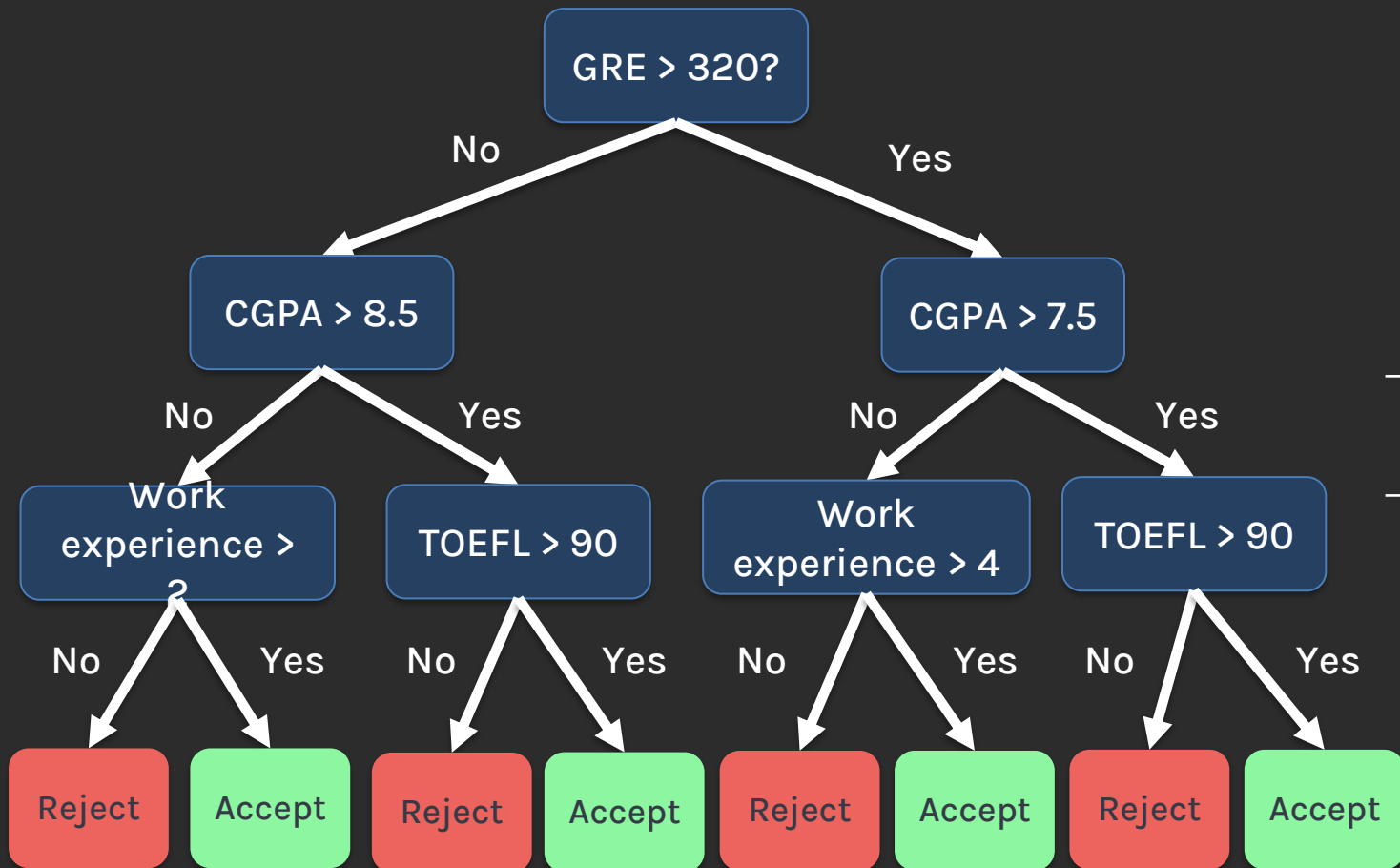
$$C(T) = \text{Error}(T) + \alpha |T|$$

The diagram illustrates the components of the cost complexity pruning formula  $C(T) = \text{Error}(T) + \alpha |T|$ . Annotations include:

- Classification Error** (orange arrow) pointing to  $\text{Error}(T)$ .
- Complexity Parameter** (purple arrow) pointing to  $\alpha$ .
- Number of leaves in the tree** (green arrow) pointing to  $|T|$ .
- Regularization term** (blue bracket and arrow) spanning the  $\alpha |T|$  portion of the equation.
- Decision tree** (white arrow) pointing to  $C(T)$ .

In other words, we add a ‘regularization’ term!

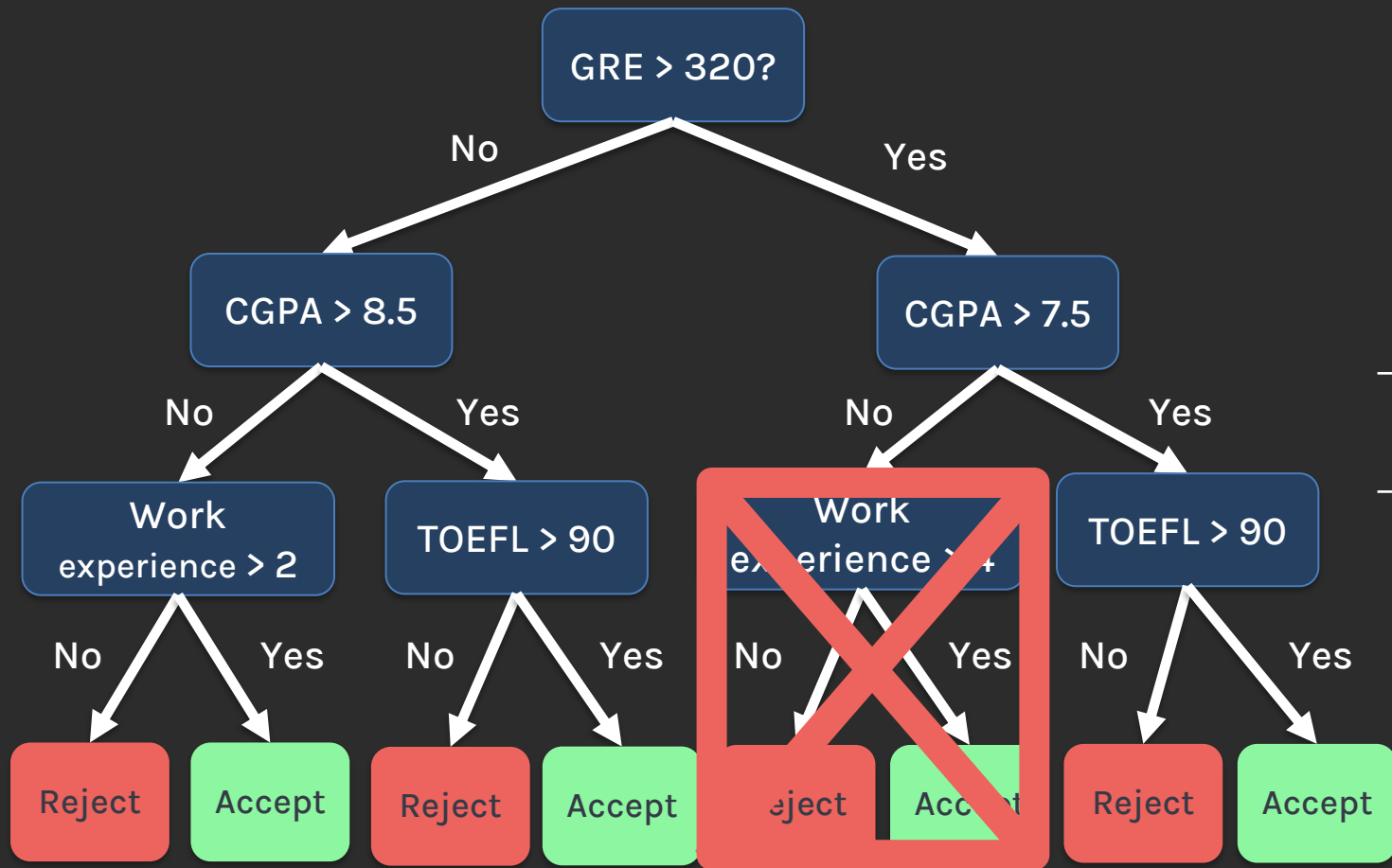
# Cost-Complexity Pruning: Example



$$\alpha = 0.2$$

Tree	Error(T)	T	Error(T) + $\alpha$  T
$T$	0.32	8	$0.32 + 0.2 * 8 = 1.92$

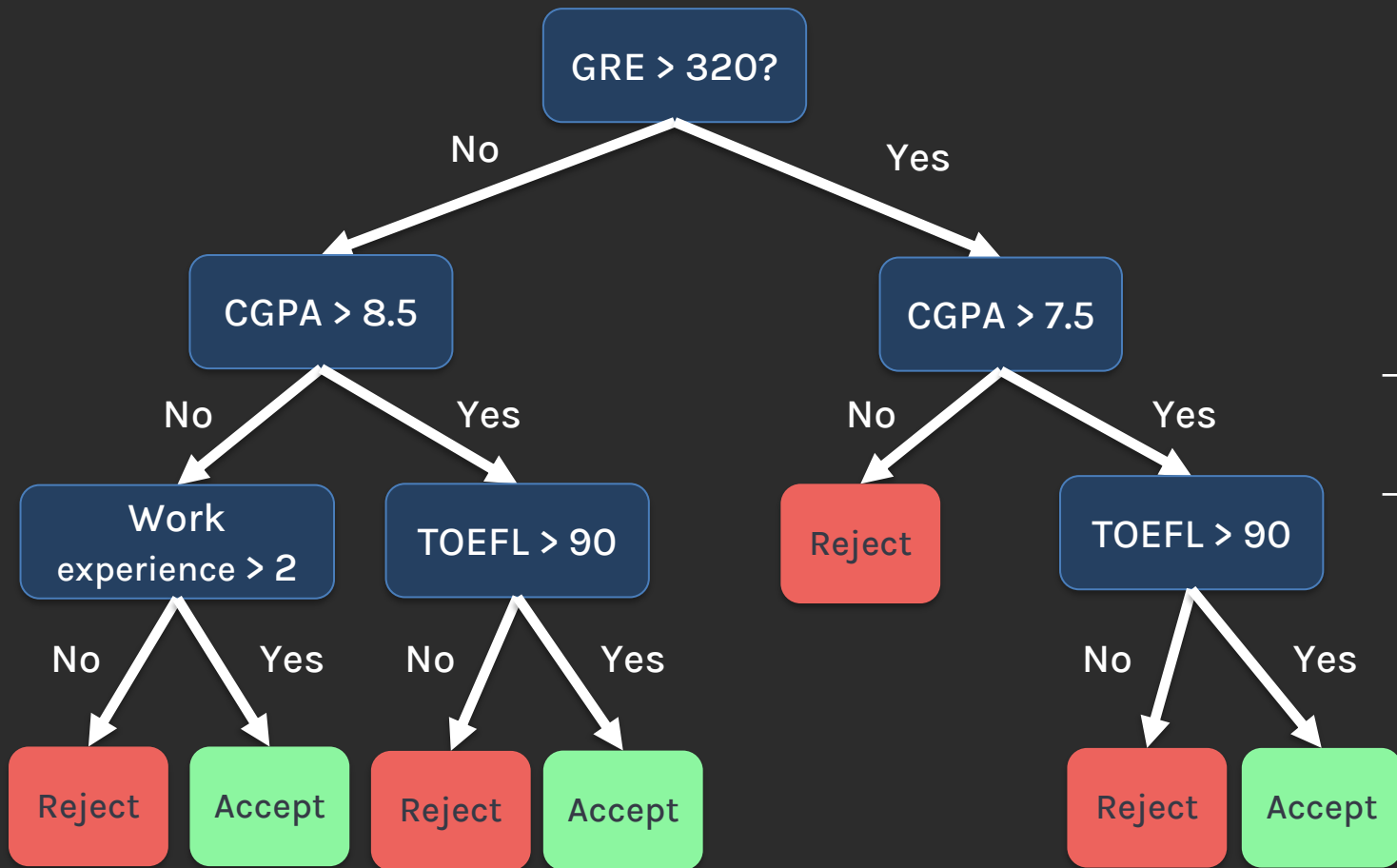
# Cost-Complexity Pruning: Example



$$\alpha = 0.2$$

Tree	Error(T)	T	Error(T) + $\alpha$  T
$T$	0.32	8	1.92

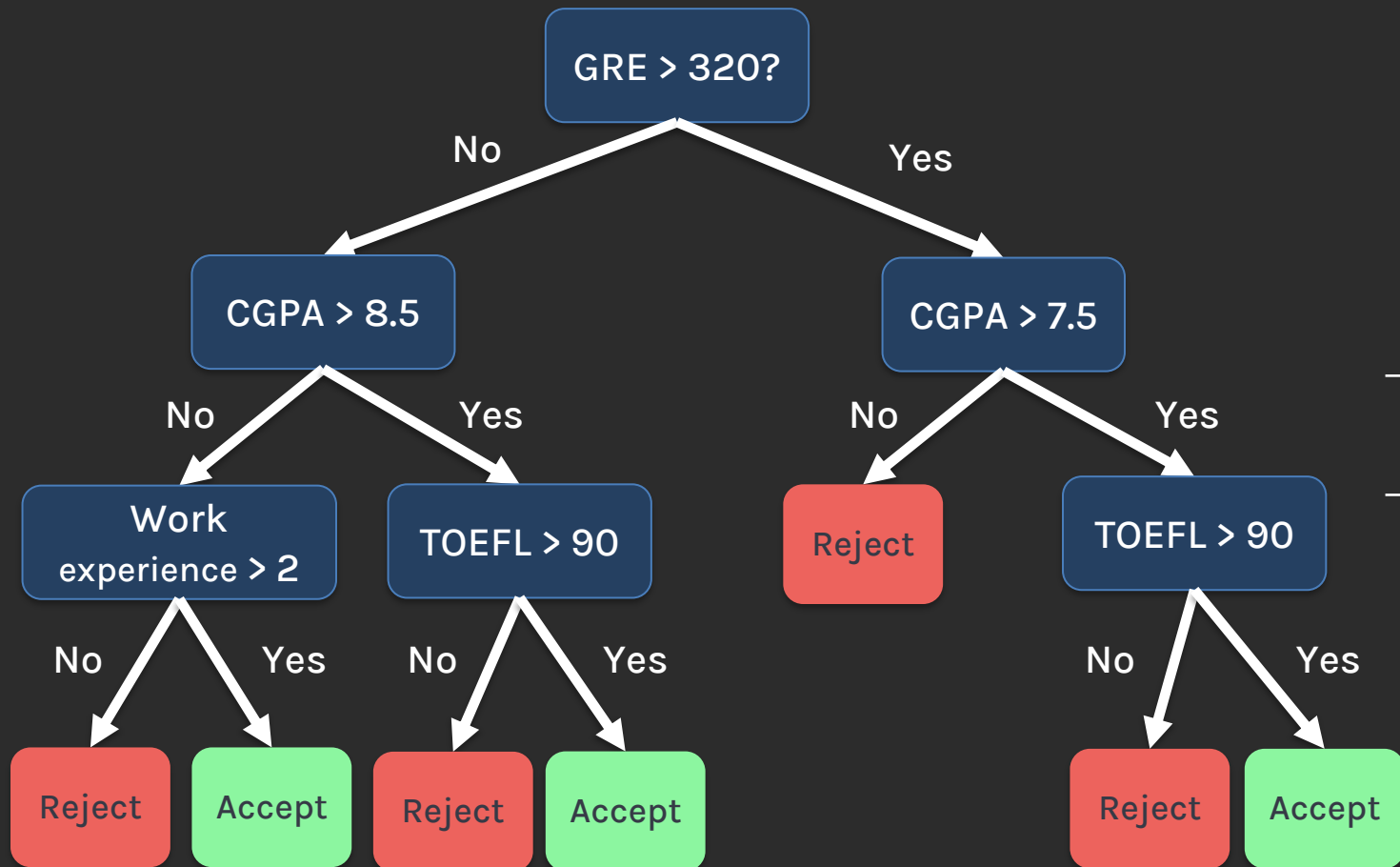
# Cost-Complexity Pruning: Example



$$\alpha = 0.2$$

Tree	Error(T)	T	Error(T) + $\alpha$  T
$T$	0.32	8	1.92
$T_{small}$	0.33	7	$0.33 + 0.2 \cdot 7 = 1.73$

# Cost-Complexity Pruning: Example



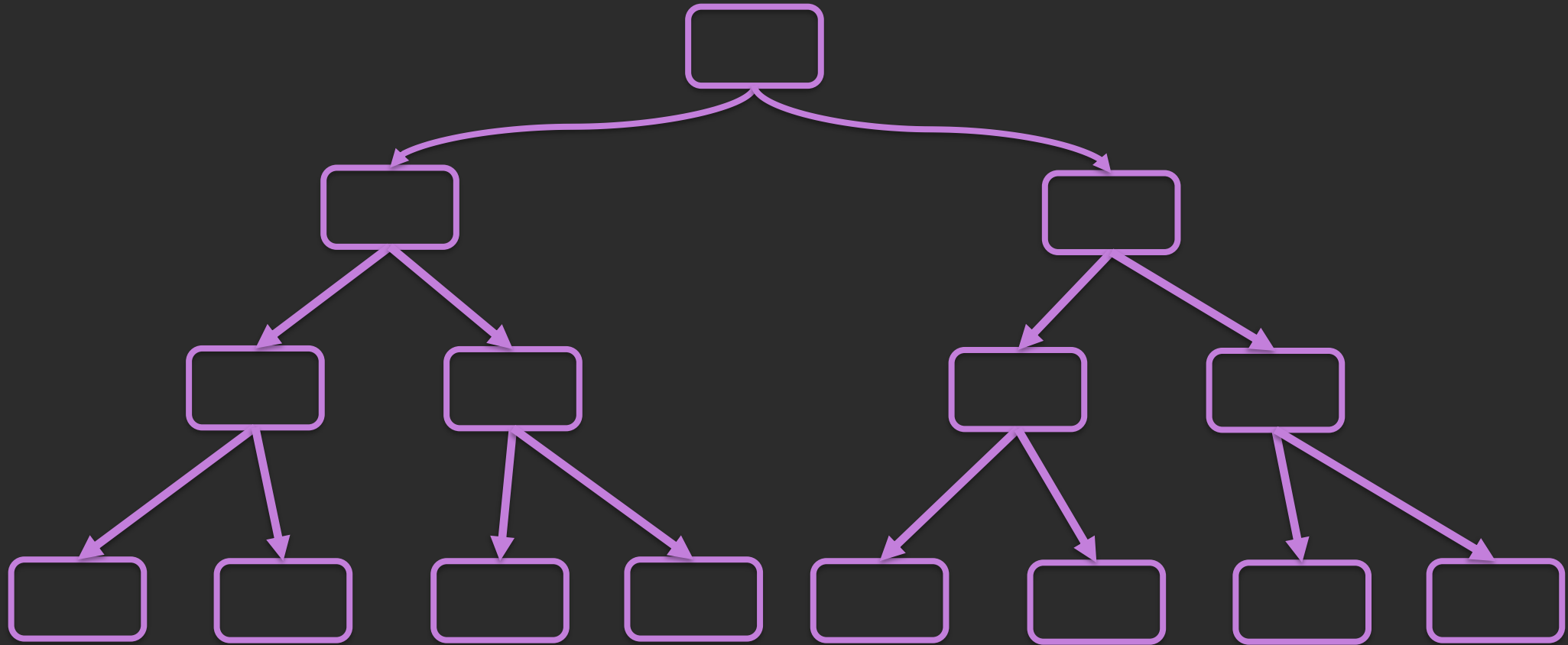
$\alpha = 0.2$

Tree	Error(T)	T	Error(T) + $\alpha$  T
$T$	0.32	8	1.92
$T_{small}$	0.33	7	1.73

The smaller tree has a larger error  $\text{Error}(T)$  but smaller complexity score  $C(T)$ .

# Pruning

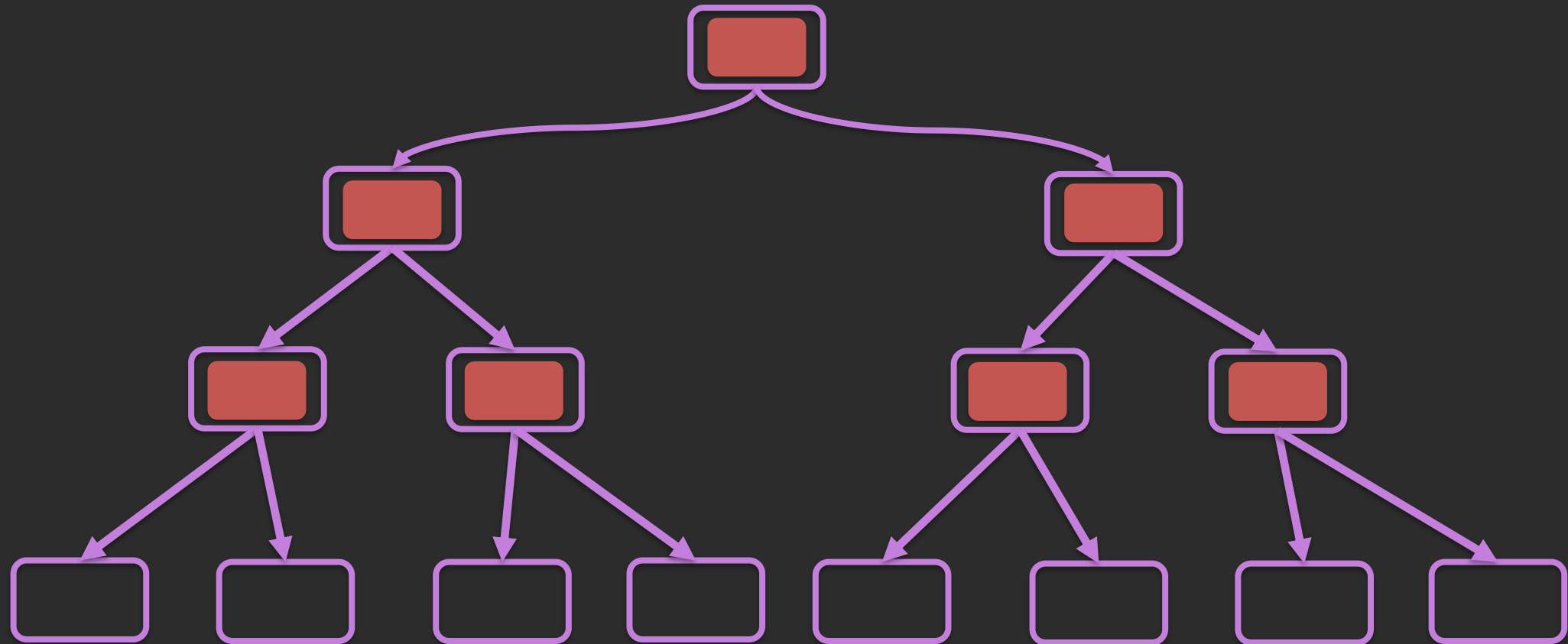
Suppose we have **a full tree**  $T_0$  as shown below.



How many possible ways are there to prune this tree?

# Pruning

There are **7** possible pruning locations, which are shown with red rectangles.

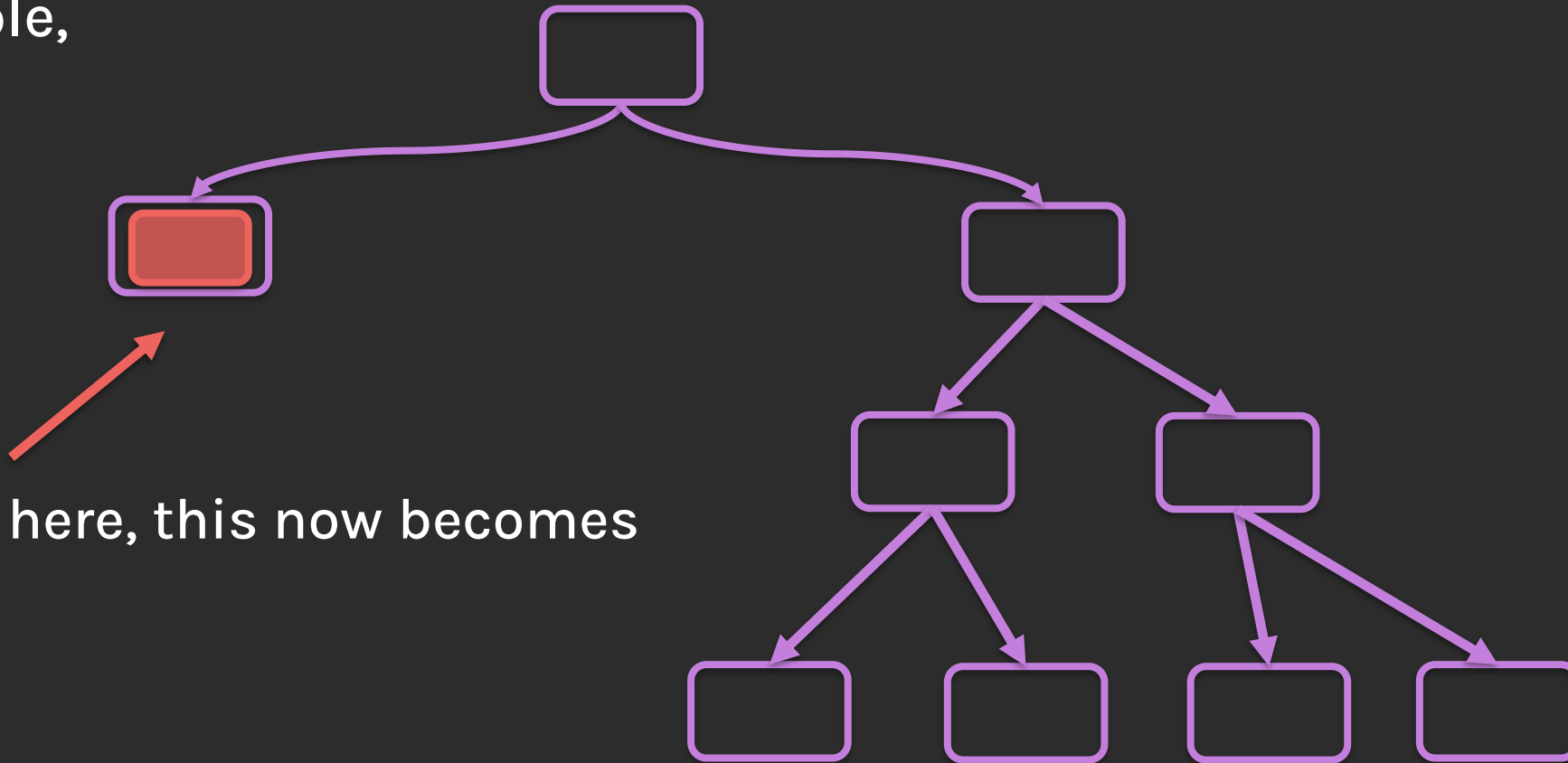




# Pruning

For each of those pruning locations, we will get 7 possible pruned trees,  $T^*$ .

For example,



If we pruned here, this now becomes a leaf node

**Question:** How do we choose the best pruned tree?

# Pruning

We will choose the one that maximizes the difference of **cost complexity score** between a full tree and a pruned tree.

Quick recap: Cost complexity score is

$$C(T) = \text{Error}(T) + \alpha |T|$$

Diagram illustrating the Cost Complexity Score formula:

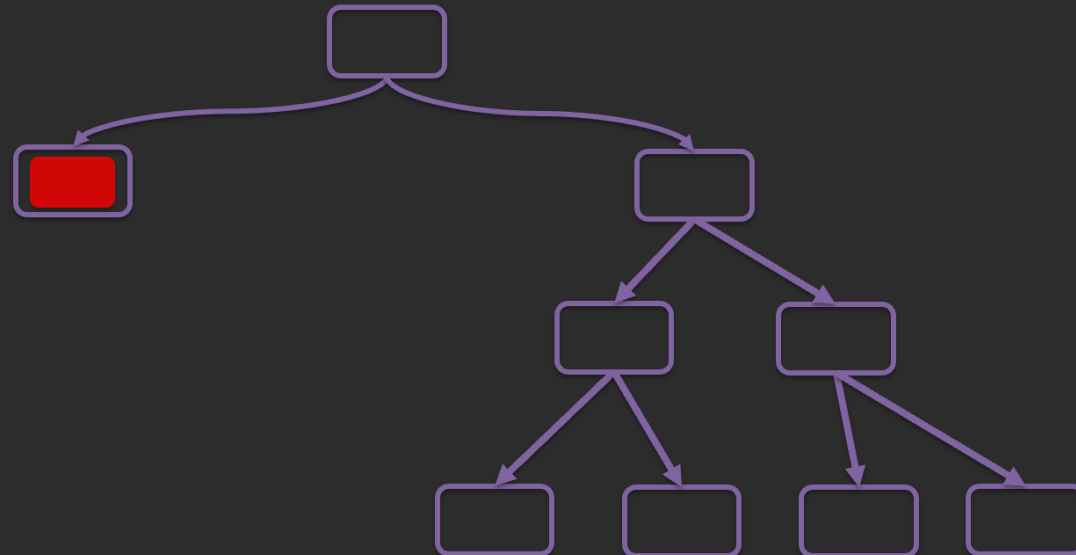
- $C(T)$ : Decision tree
- $\text{Error}(T)$ : Classification Error
- $\alpha$ : Complexity Parameter
- $|T|$ : Number of leaves in the tree

# Pruning

Our goal is to maximize  $\mathcal{C}(T) - \mathcal{C}(T^*)$

1.  $\operatorname{argmax}_{T^*} [\mathcal{C}(T) - \mathcal{C}(T^*)]$
2.  $\operatorname{argmax}_{T^*} [E(T) - E(T^*) + \alpha|T| - \alpha|T^*|]$
3.  $\operatorname{argmax}_{T^*} \left[ \frac{E(T) - E(T^*)}{\alpha|T^*| - \alpha|T|} - 1 \right]$  (divide by  $\alpha|T^*| - \alpha|T|$ )
4.  $\operatorname{argmax}_{T^*} \left[ \frac{E(T) - E(T^*)}{\alpha|T^*| - \alpha|T|} \right]$
5.  $\operatorname{argmin}_{T^*} \left[ \frac{E(T) - E(T^*)}{\alpha|T| - \alpha|T^*|} \right]$

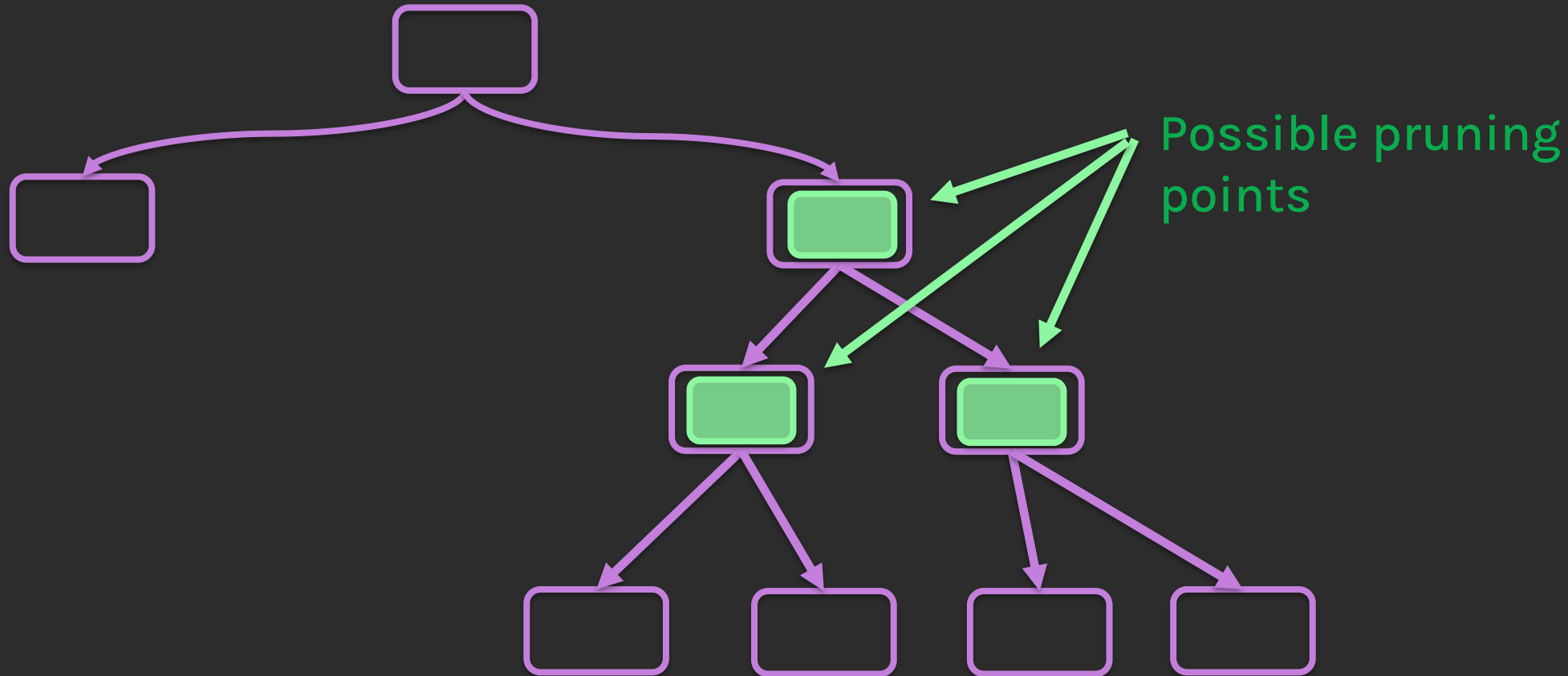
$T^*$  is the pruned tree where  $T$  is the unpruned tree



This will result in a subtree,  $T^{(1)}$

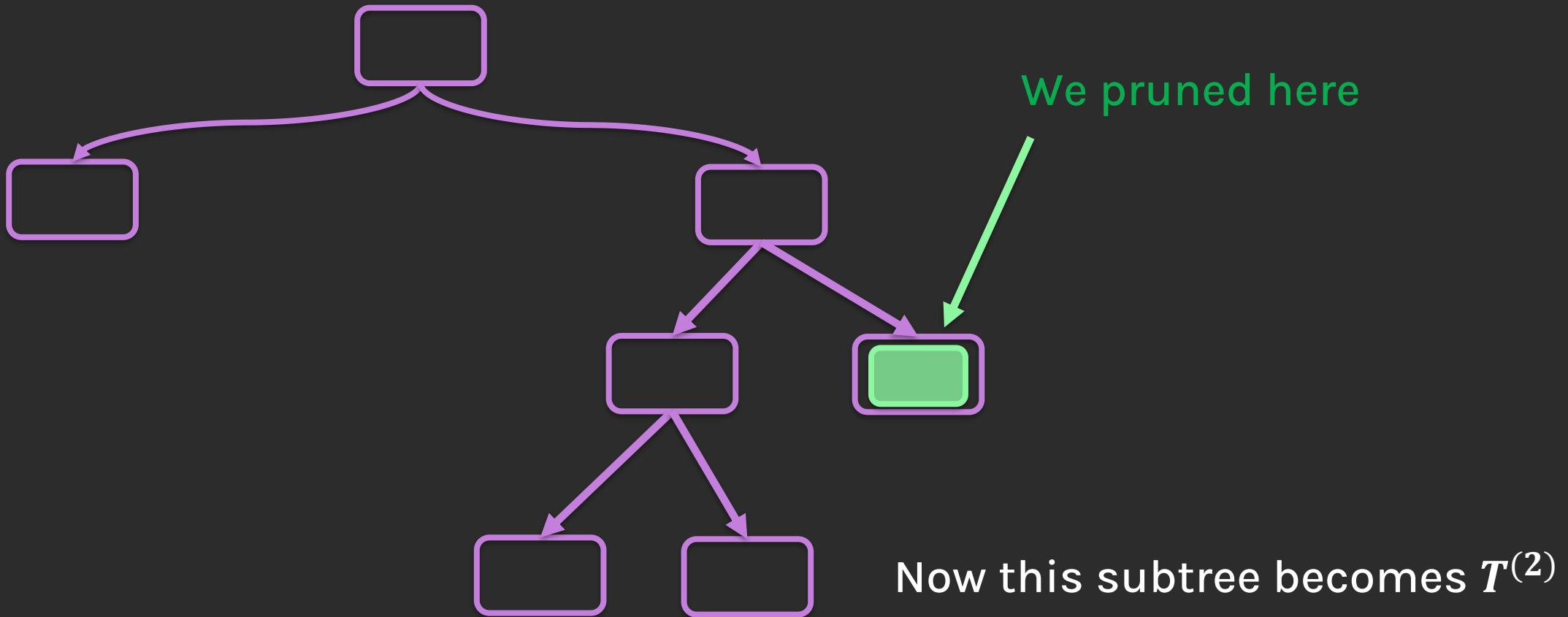
# Pruning

Now, we again consider all possible pruning from  $T^{(1)}$ .



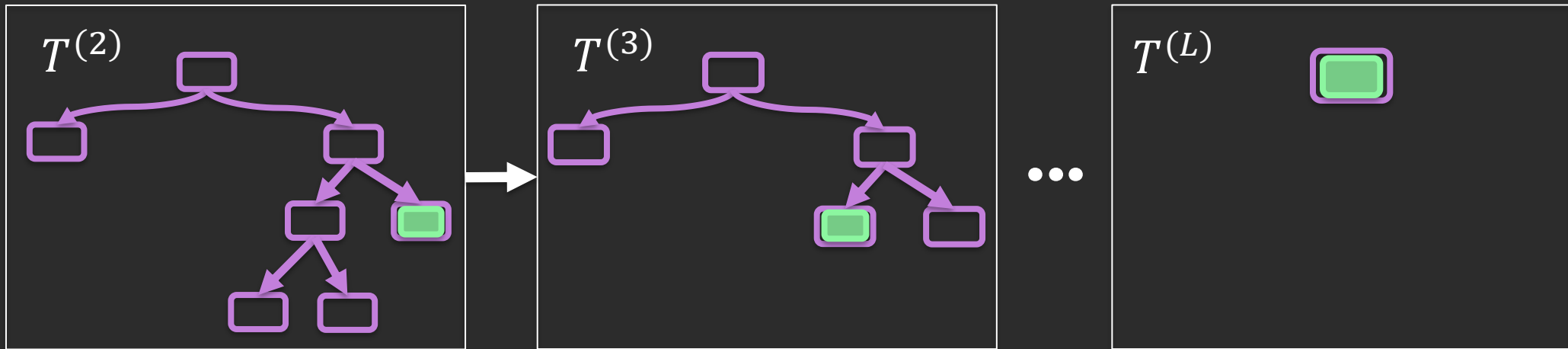
# Pruning

We again minimize the ratio of the difference of the errors **over the difference in the complexity**

$$\underset{T^*}{\operatorname{argmin}} \left[ \frac{E(T^{(1)}) - E(T^*)}{\alpha|T^{(1)}| - \alpha|T^*|} \right]$$


# Pruning

We iterate this pruning process to obtain  $T^{(2)}, T^{(3)}, \dots, T^{(L)}$  where  $T^{(L)}$  is the tree containing just the root of  $T^{(0)}$ .



We select the optimal tree  $T^{(i)}$  by cross validation.

This is the  $T^*$  given  $\alpha$ . Finally, we choose **the optimal**  $\alpha$  by cross validation!

# Summary

What is the main drawback of pre-defining stopping criteria for decision tree growth?

Pre-defined stopping criteria may lead to trees that are either too simple (underfitting) or too complex (overfitting) since it's challenging to determine the optimal stopping point beforehand.

Explain the alternative approach to using stopping conditions for decision tree construction.

Instead of stopping conditions, an alternative approach is to grow a large tree first, then prune it back to find a balanced structure, offering flexibility for optimal complexity.

What is the core concept behind pruning in decision trees?

Pruning simplifies a decision tree by removing branches with minimal impact on overall accuracy, enhancing the tree's ability to generalize.



# Summary

Describe the formula used for calculating the cost complexity of a decision tree.

Cost complexity ( $C(T)$ ) is calculated as:  $C(T) = \text{Error}(T) + \alpha|T|$ , where  $\text{Error}(T)$  is the classification error,  $\alpha$  is the complexity parameter, and  $|T|$  is the number of leaves.

In cost-complexity pruning, how does the complexity parameter ( $\alpha$ ) influence the trade-off between tree size and error?

The complexity parameter ( $\alpha$ ) adjusts the penalty for tree size, with higher values favoring smaller trees with potentially higher error and lower values favoring more complex trees with potentially lower error.

How do you determine the best pruned tree from a set of candidate trees generated by pruning?

The best pruned tree is selected by maximizing the difference in cost complexity scores, choosing the tree with the most significant reduction in complexity without a major increase in error.



Thank you

