

# Decision Trees

CS771: Introduction to Machine Learning

Purushottam Kar

# Announcements

2

Please form groups of 5 by Wednesday – will give you an online form to specify your group details – please do not mail group details to me

Please activate your Piazza account if not done so

Heads-up: Quiz 1 is coming up next week August 14 – do not miss



# Recap of Last Lecture

3

Notion of decision boundary – points where the classifiers are confused about what to predict - all classifiers have a decision boundary whether LwP, NN, trees, deep nets

Hyperplane classifiers – model vector, bias and their role

Metric learning and how it can help LwP and NN improve

NN variants – kNN, rNN (not to be confused with RNN)

Hyperparameter tuning via held out validation/cross validation



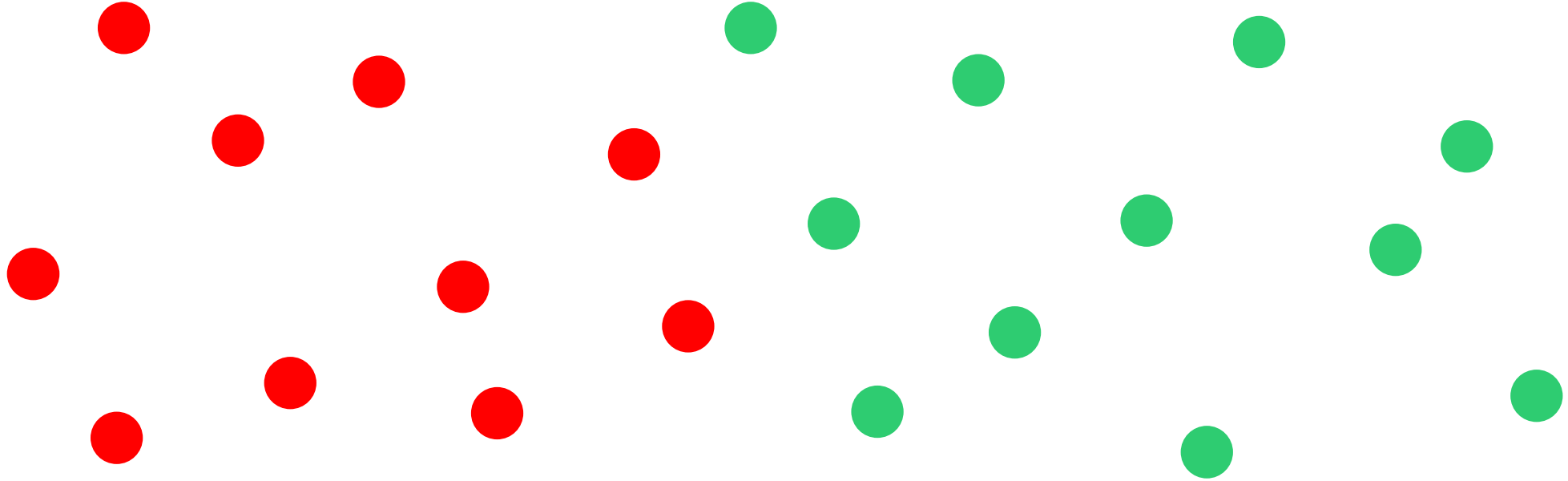
# Regression with rNN

4



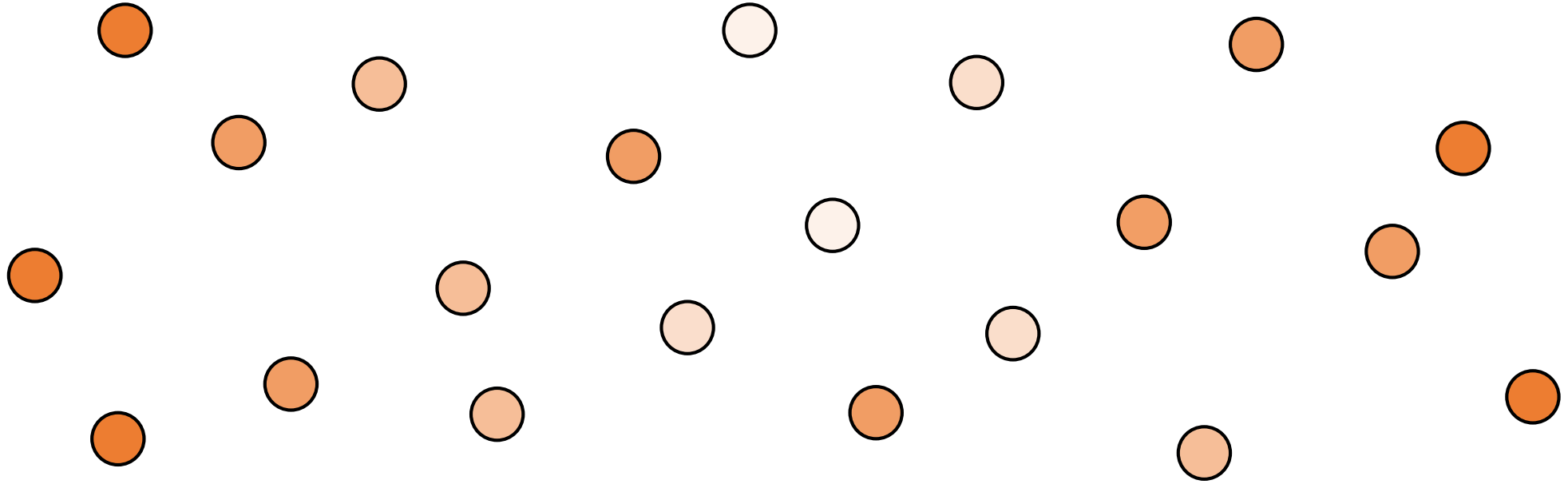
# Regression with rNN

4



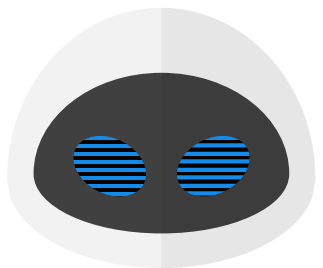
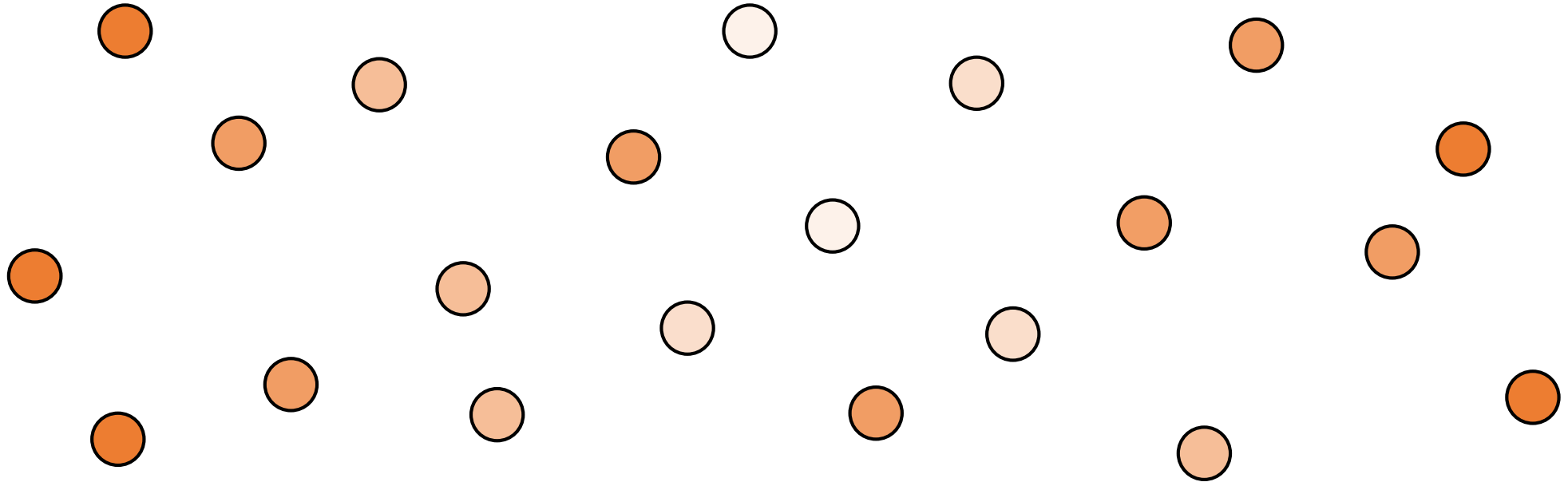
# Regression with rNN

4



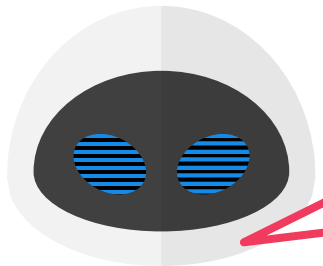
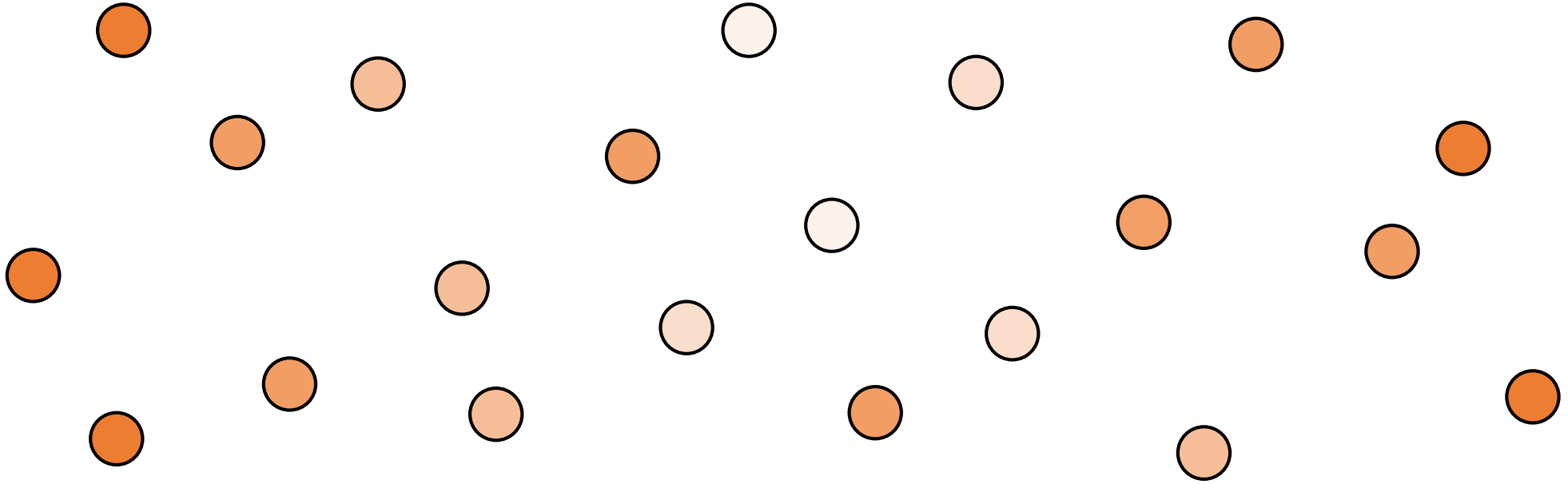
# Regression with rNN

4



# Regression with rNN

4



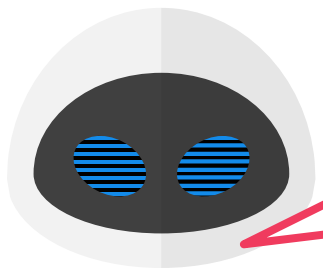
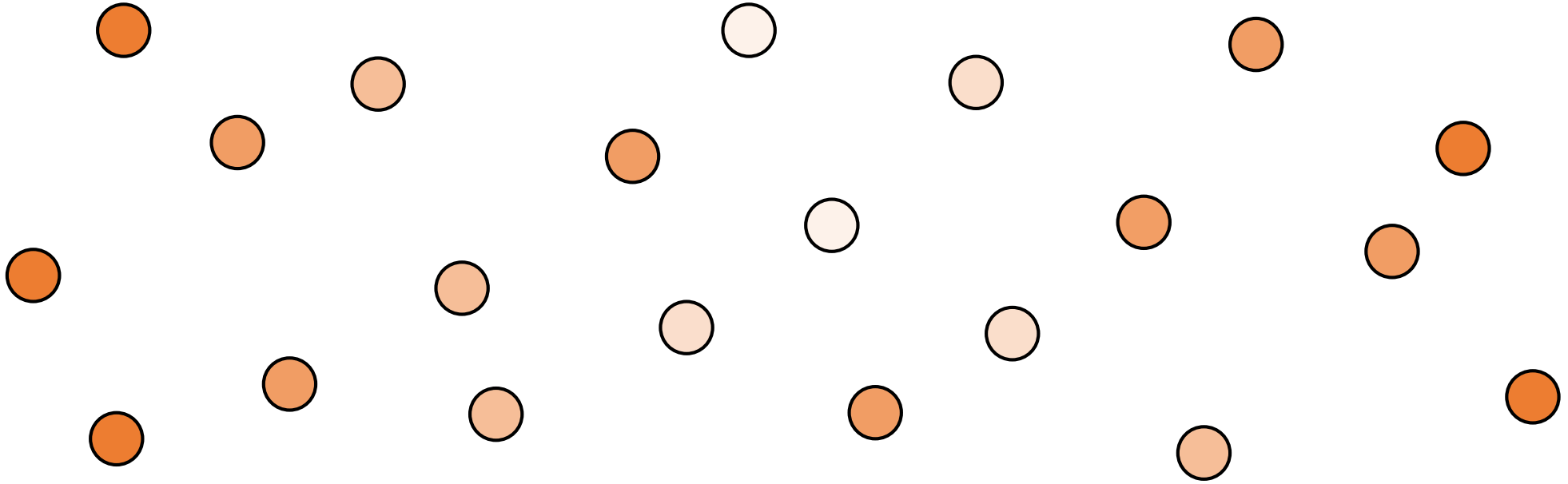
In regression, each point needs to be given a real-valued score instead of a label like spam/non-spam





# Regression with rNN

4

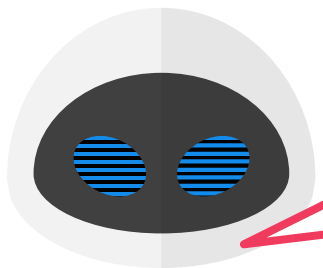
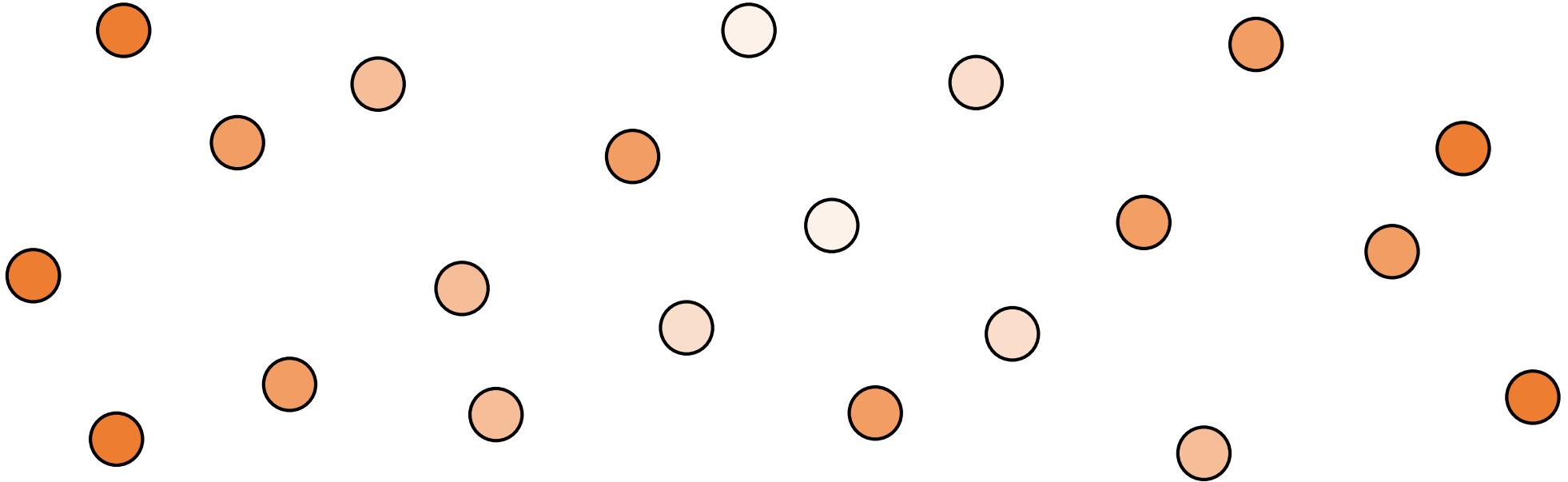


In regression, each point needs to be given a real-valued score instead of a label like spam/non-spam



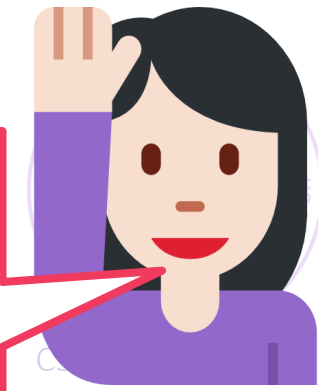
# Regression with rNN

4



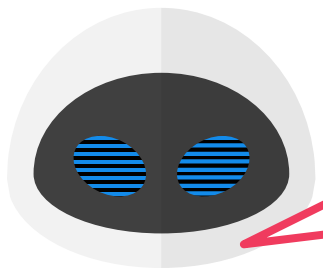
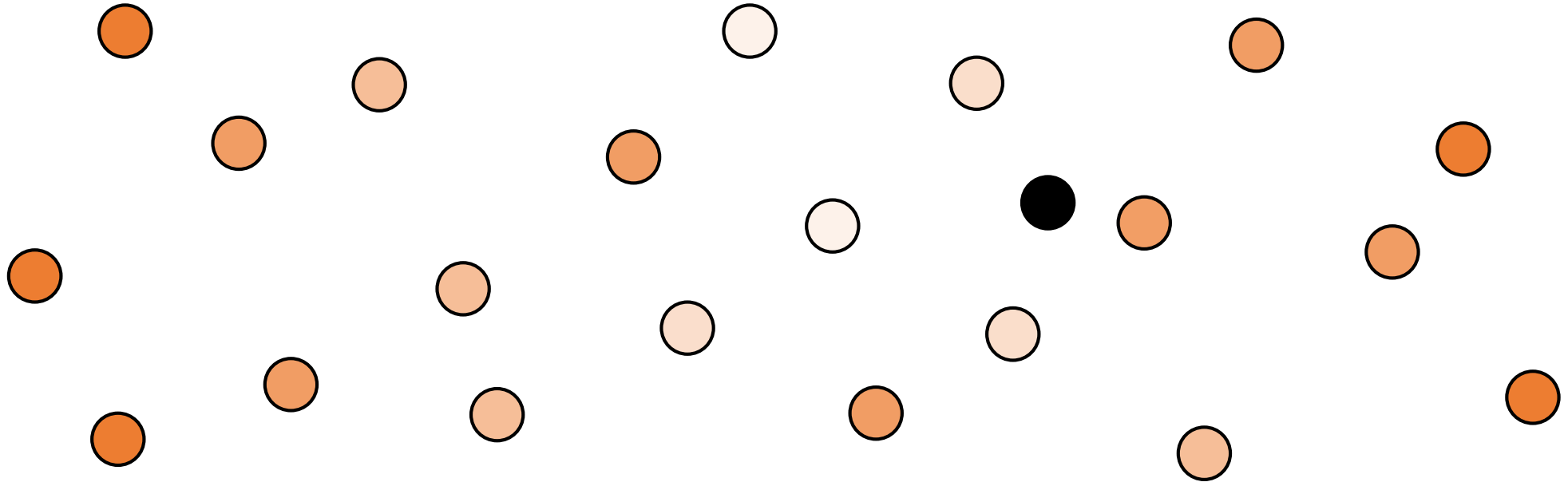
In regression, each point needs to be given a real-valued score instead of a label like spam/non-spam

E.g. when data point is a student and we wish to predict their marks



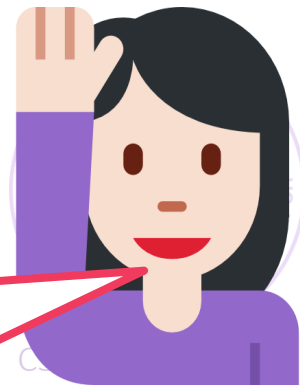
# Regression with rNN

4



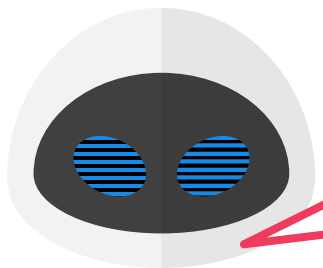
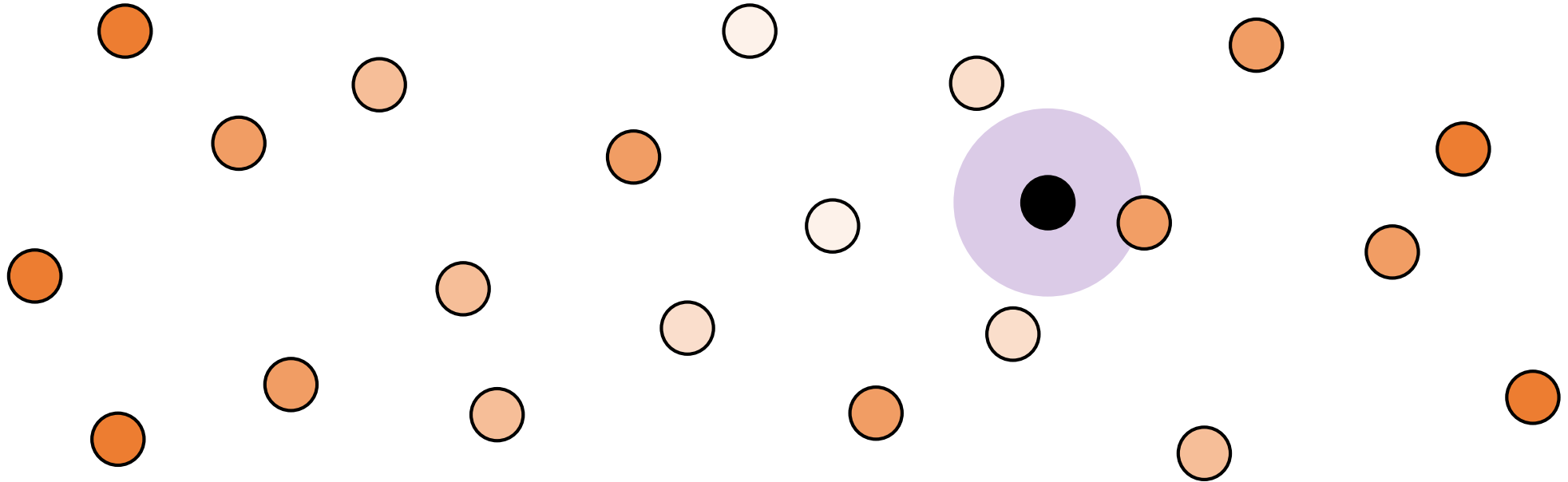
In regression, each point needs to be given a real-valued score instead of a label like spam/non-spam

E.g. when data point is a student and we wish to predict their marks



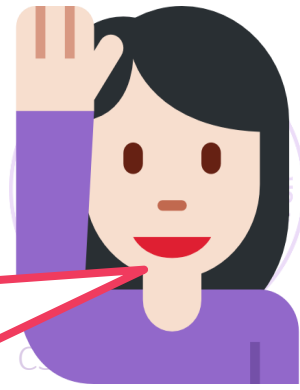
# Regression with rNN

4



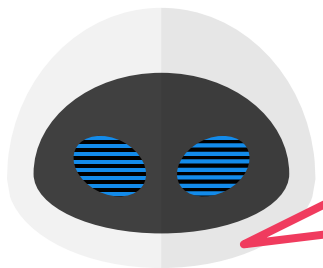
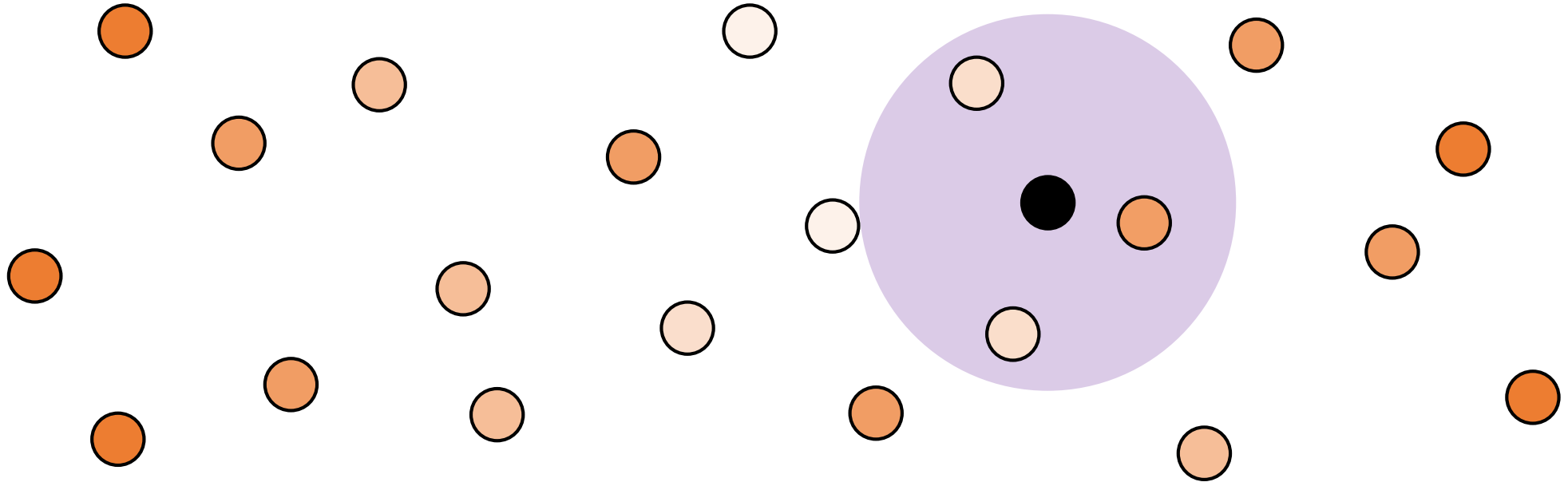
In regression, each point needs to be given a real-valued score instead of a label like spam/non-spam

E.g. when data point is a student and we wish to predict their marks



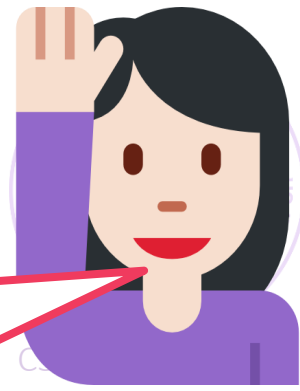
# Regression with rNN

4



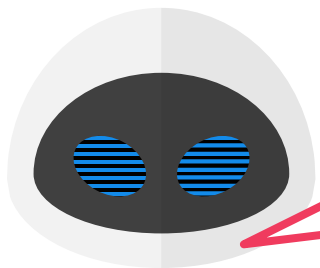
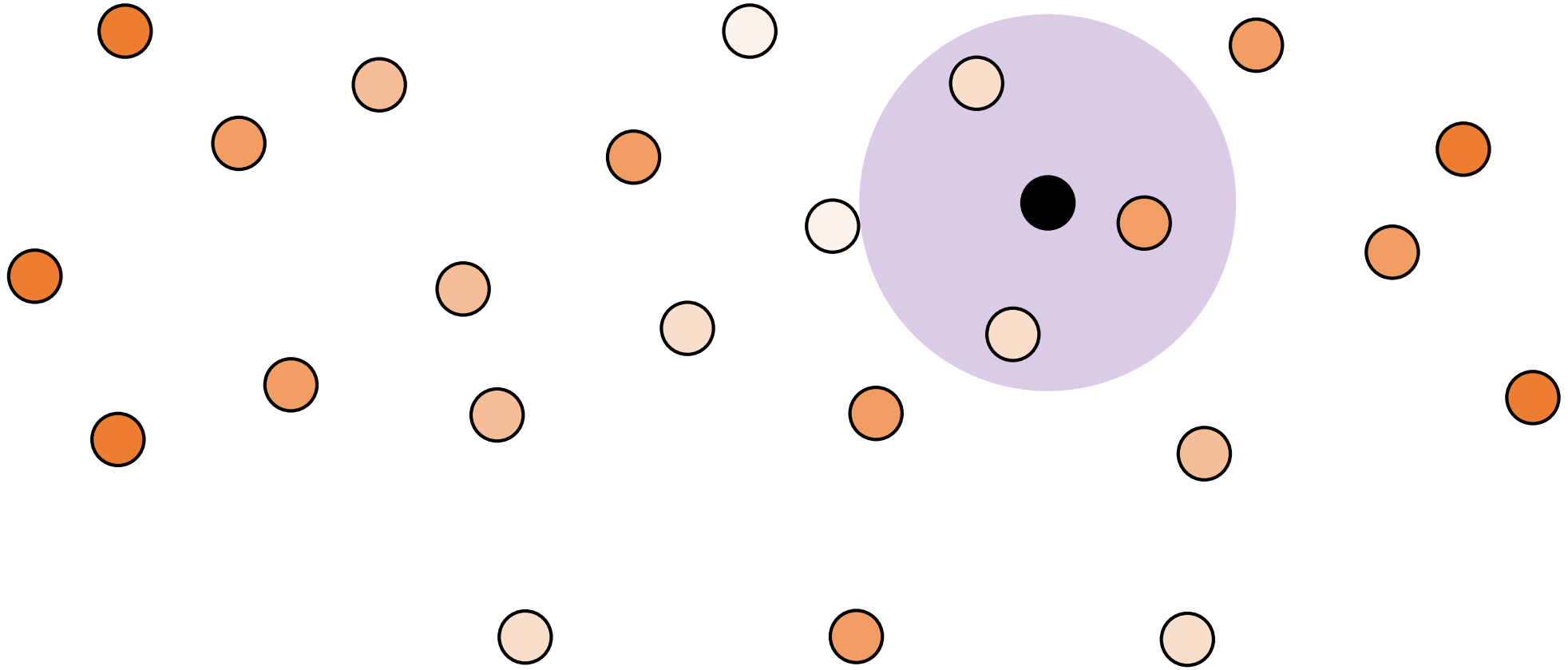
In regression, each point needs to be given a real-valued score instead of a label like spam/non-spam

E.g. when data point is a student and we wish to predict their marks



# Regression with rNN

4



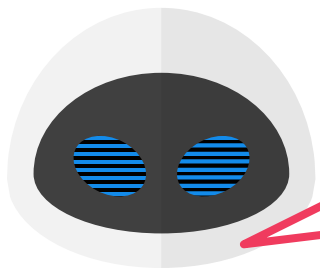
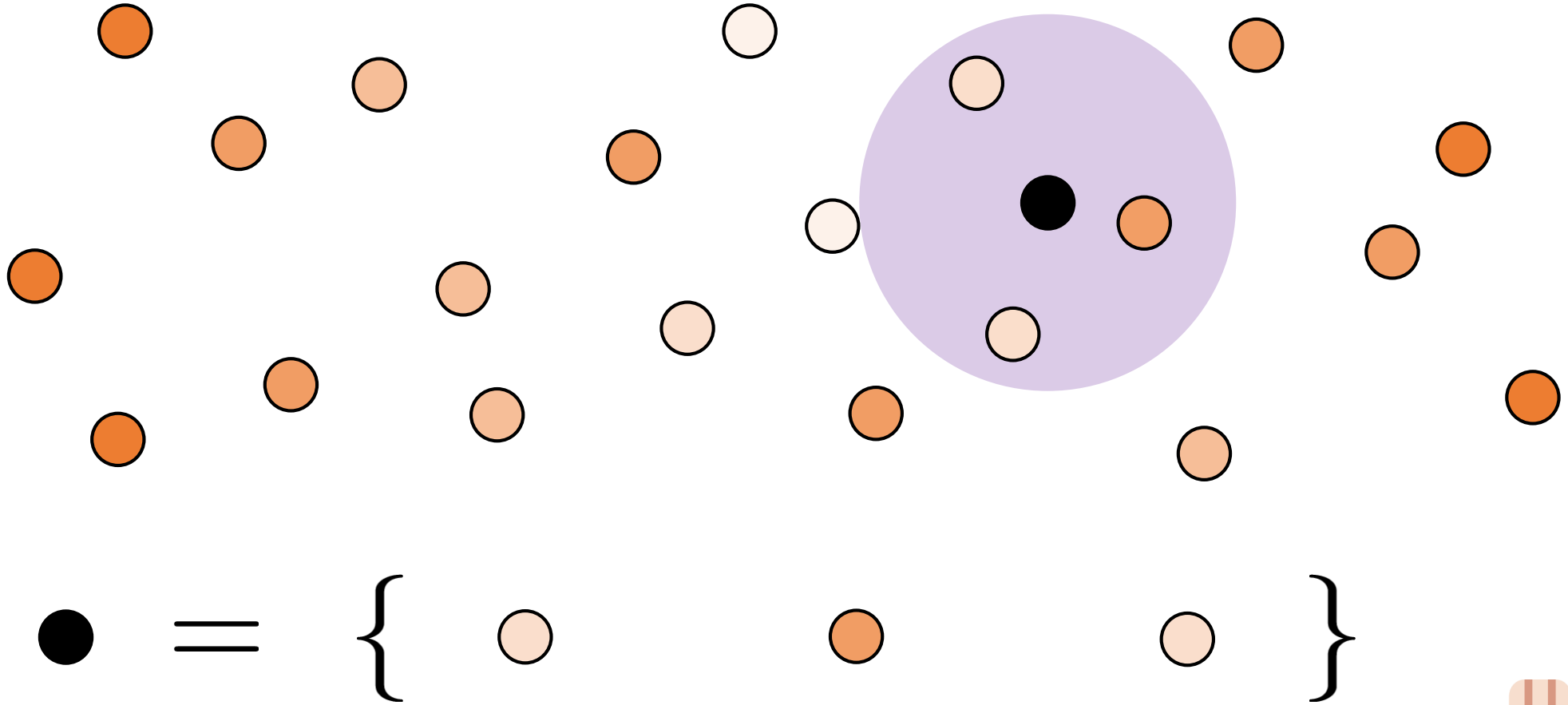
In regression, each point needs to be given a real-valued score instead of a label like spam/non-spam

E.g. when data point is a student and we wish to predict their marks



# Regression with rNN

4

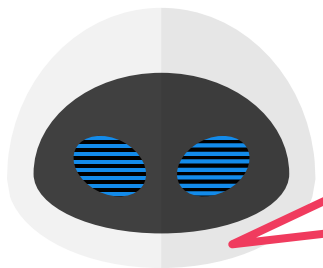
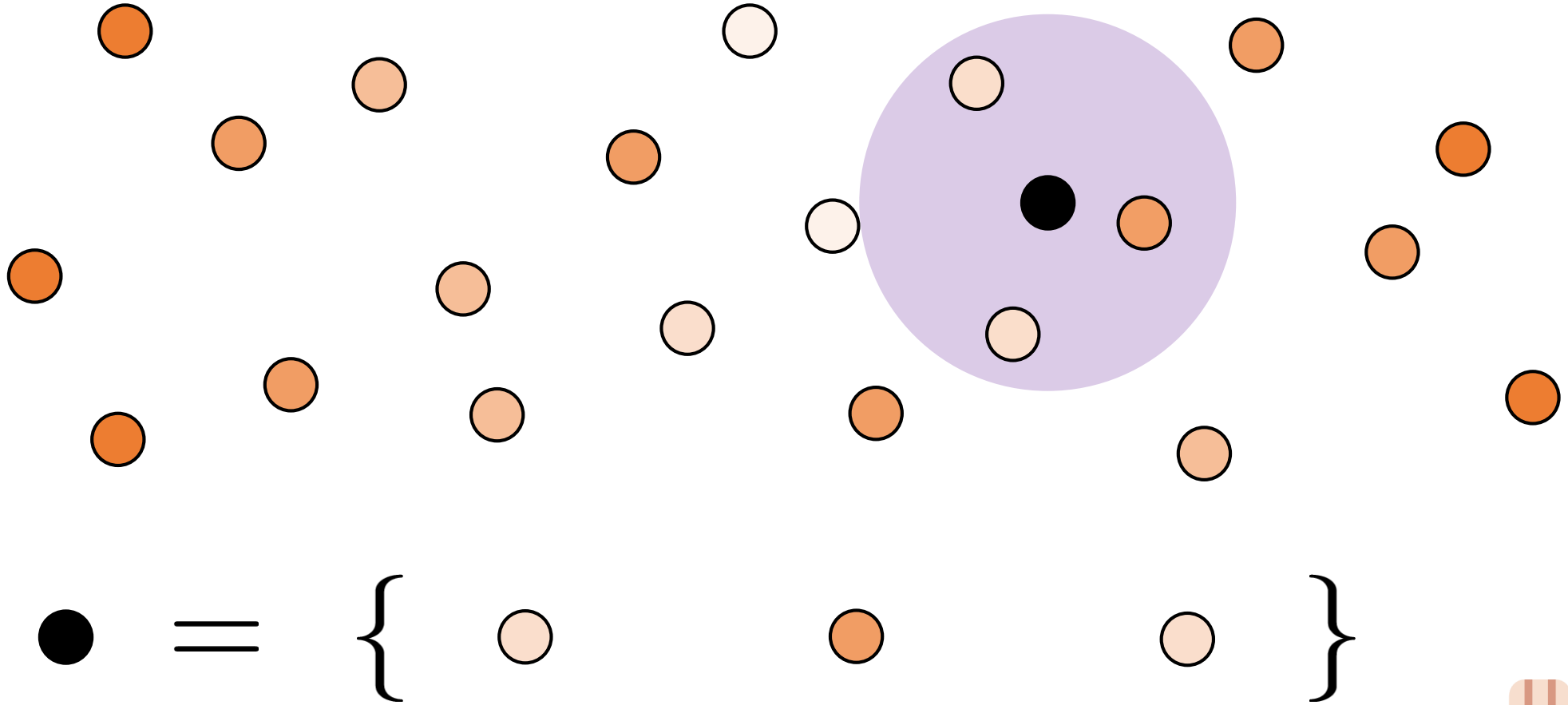


In regression, each point needs to be given a real-valued score instead of a label like spam/non-spam

E.g. when data point is a student and we wish to predict their marks

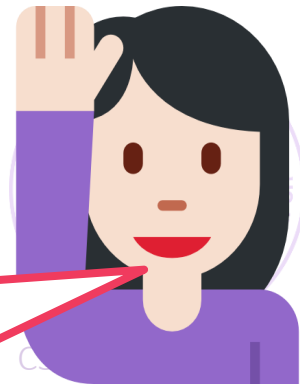


# Regression with rNN



Once you have found the neighbours,  
you can average their scores to  
predict a score for the test data point

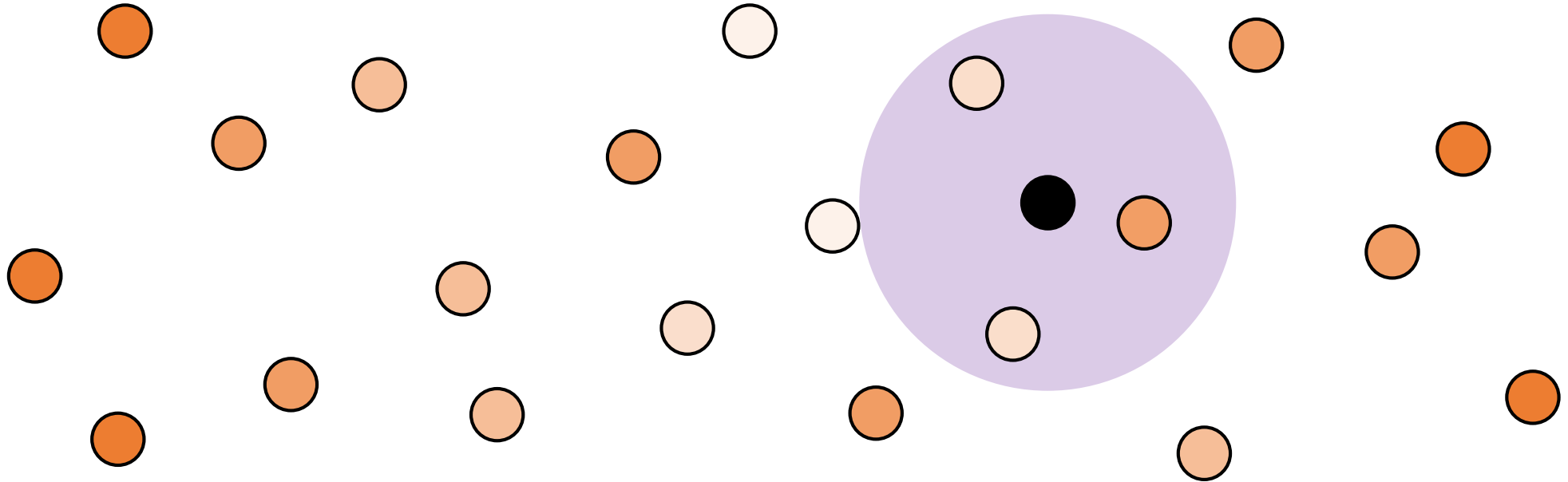
E.g. when data point is a student and we wish to predict their marks



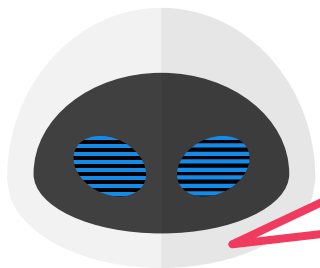


# Regression with rNN

4



$$\bullet = \left\{ \frac{1}{3} \text{ (light orange dot)} + \frac{1}{3} \text{ (dark orange dot)} + \frac{1}{3} \text{ (light orange dot)} \right\}$$



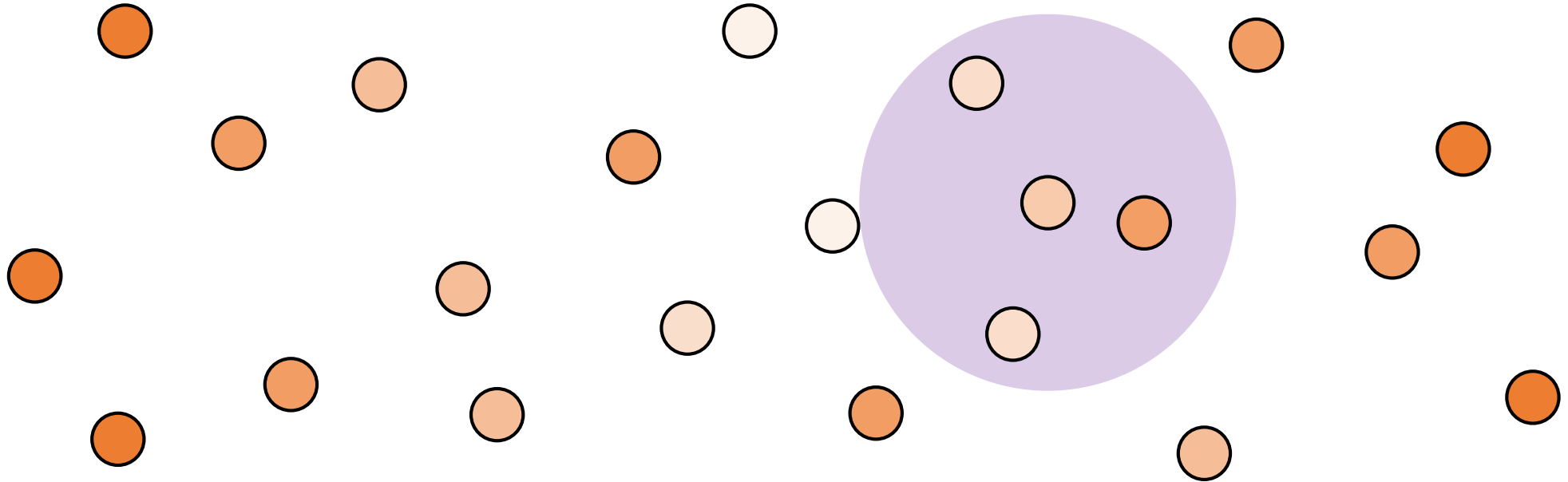
Once you have found the neighbours, you can average their scores to predict a score for the test data point

E.g. when data point is a student and we wish to predict their marks

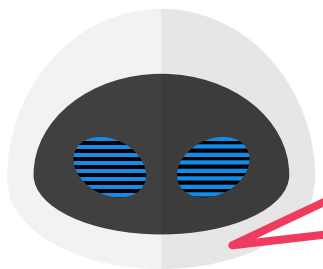


# Regression with rNN

4



$$\text{light orange circle} = \left\{ \frac{1}{3} \text{light orange circle} + \frac{1}{3} \text{orange circle} + \frac{1}{3} \text{light orange circle} \right\}$$



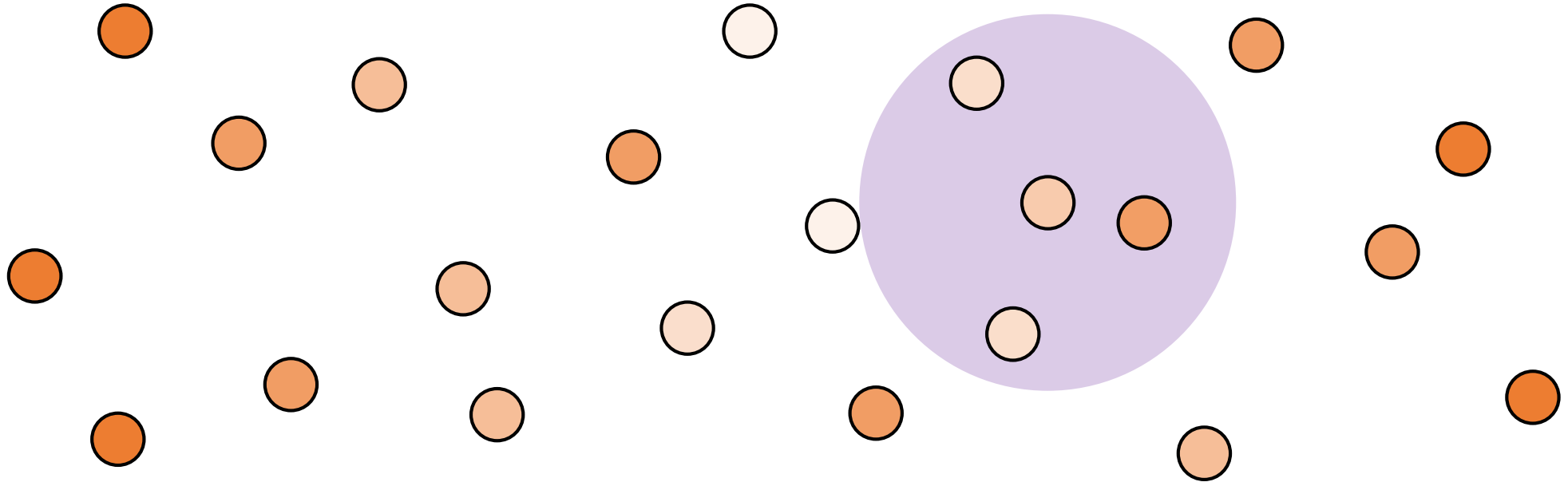
Once you have found the neighbours, you can average their scores to predict a score for the test data point

E.g. when data point is a student and we wish to predict their marks

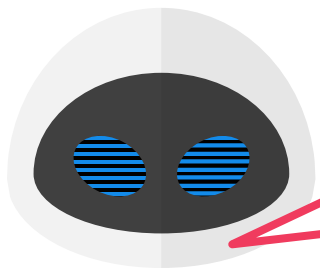


# Regression with Weighted rNN

4



$$\text{light orange circle} = \left\{ \frac{1}{3} \text{light orange circle} + \frac{1}{3} \text{orange circle} + \frac{1}{3} \text{light orange circle} \right\}$$



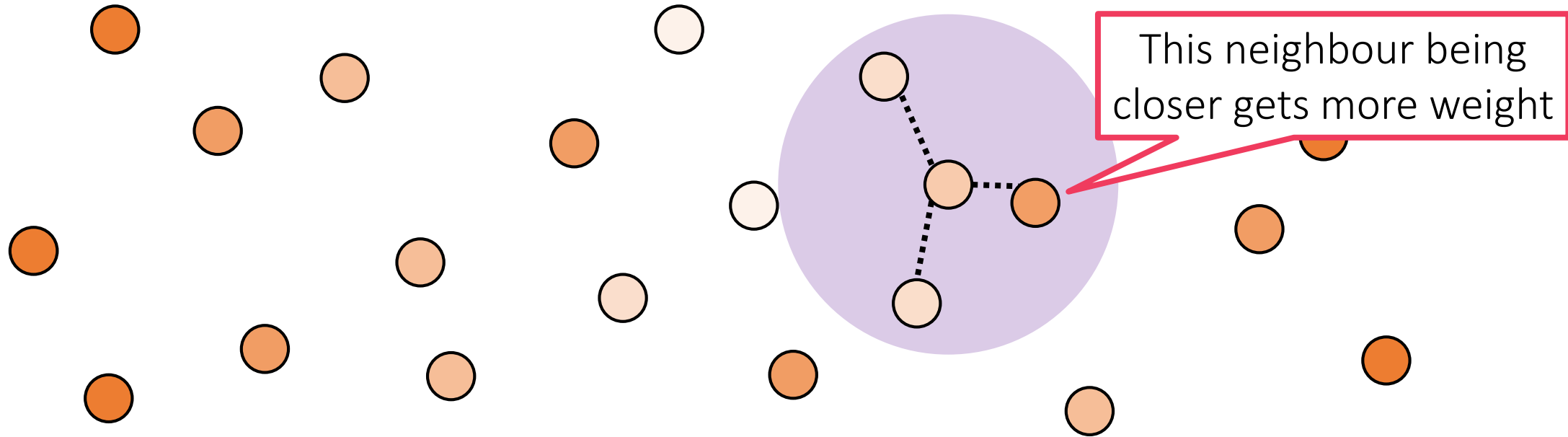
Once you have found the neighbours, you can average their scores to predict a score for the test data point

E.g. when data point is a student and we wish to predict their marks

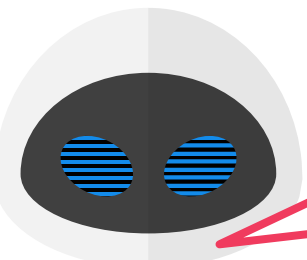


# Regression with Weighted rNN


4



$$\text{light orange circle} = \left\{ \frac{1}{3} \text{light orange circle} + \frac{1}{3} \text{orange circle} + \frac{1}{3} \text{light orange circle} \right\}$$



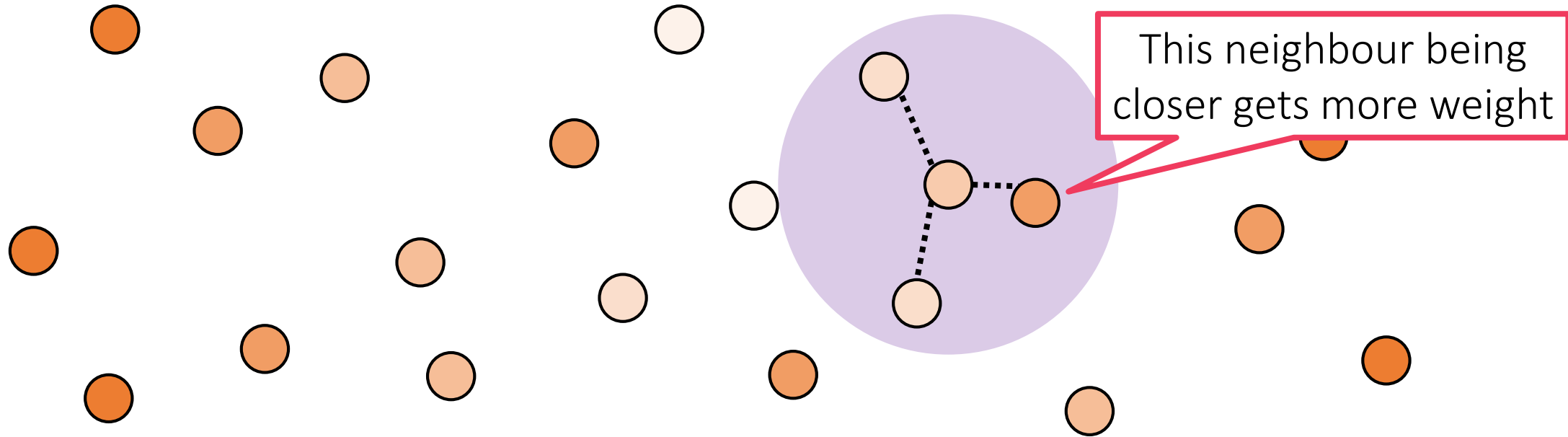
Or else, you may give the score of closer neighbours more weight and those of far neighbours less weight



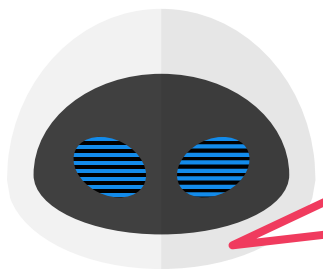
E.g. when data point is a student and we wish to predict their marks

# Regression with Weighted rNN

4



$$\text{light orange circle} = \left\{ \frac{1}{5} \text{light orange circle} + \frac{3}{5} \text{orange circle} + \frac{1}{5} \text{light orange circle} \right\}$$



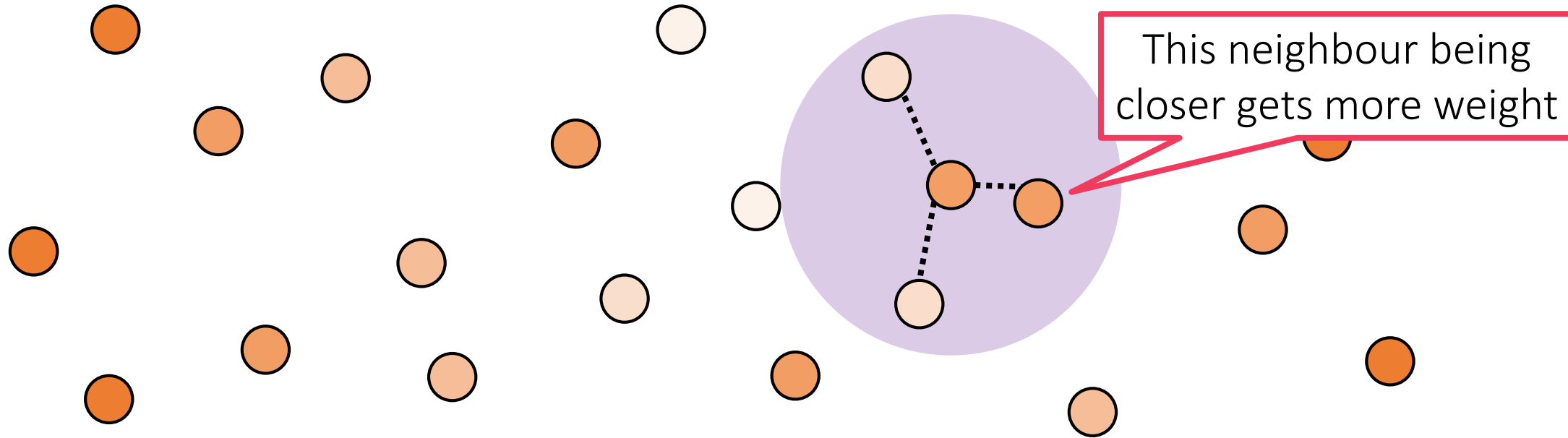
Or else, you may give the score of closer neighbours more weight and those of far neighbours less weight

E.g. when data point is a student and we wish to predict their marks

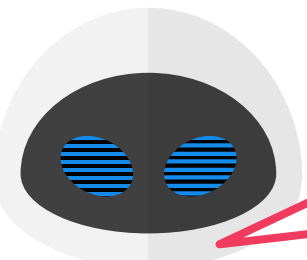


# Regression with Weighted rNN

4




$$\text{orange circle} = \left\{ \frac{1}{5} \text{light orange circle} + \frac{3}{5} \text{orange circle} + \frac{1}{5} \text{light orange circle} \right\}$$



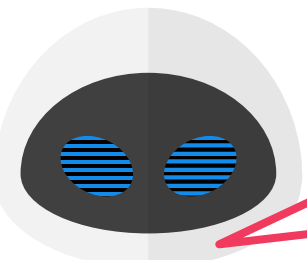
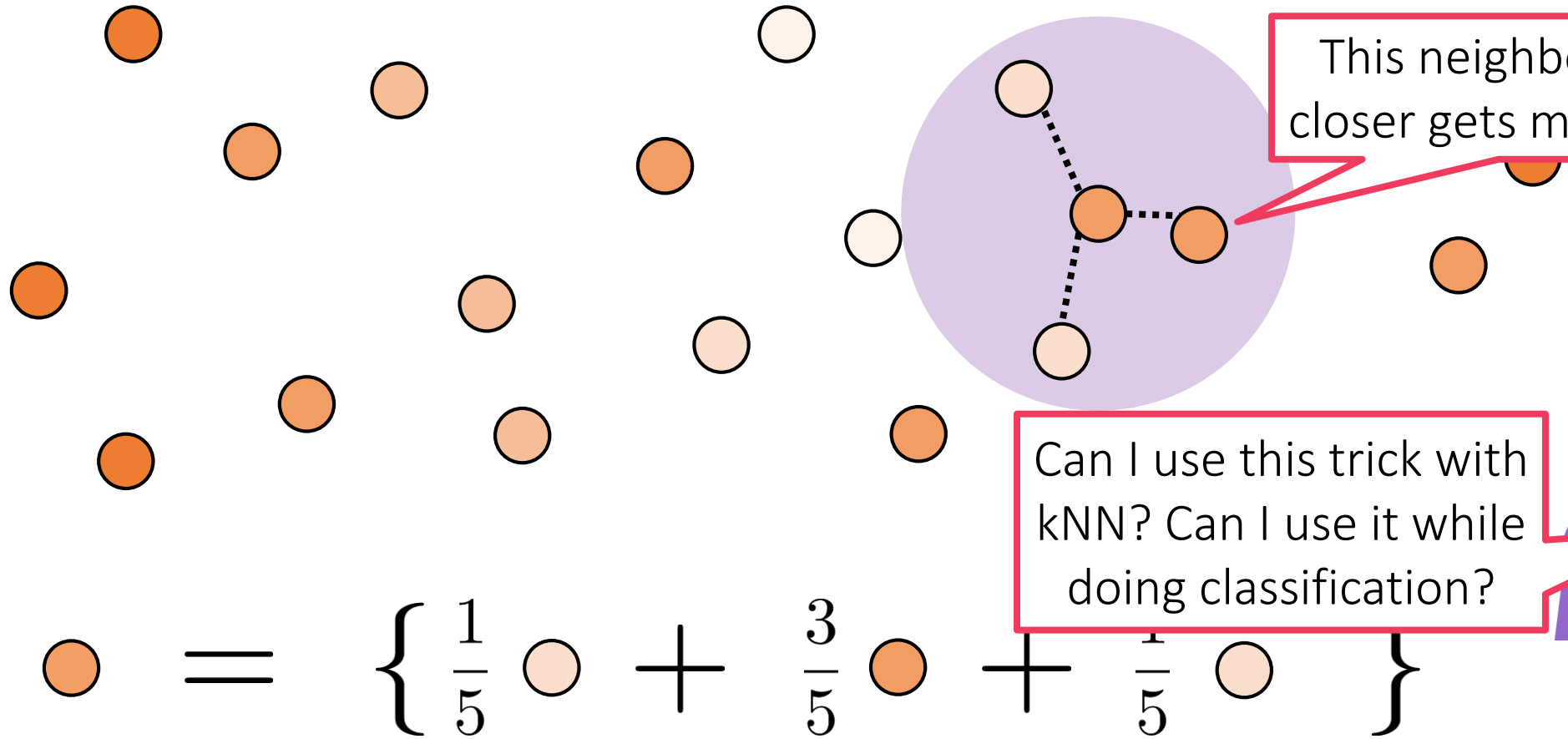
Or else, you may give the score of closer neighbours more weight and those of far neighbours less weight

E.g. when data point is a student and we wish to predict their marks

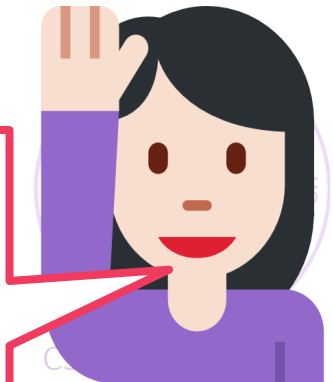


# Regression with Weighted rNN

4



Yes, and yes. Need to be a bit careful while doing weighted classification since adding labels makes no sense



E.g. when data point is a student and we wish to predict their marks

# Classification with Weighted rNN

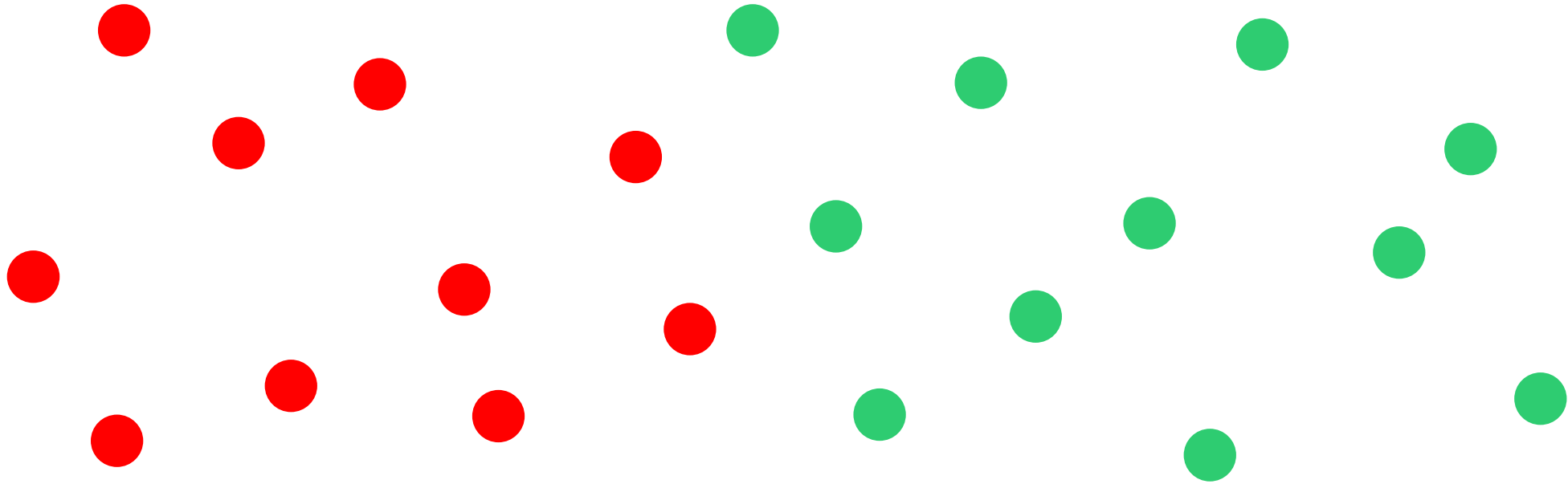
26





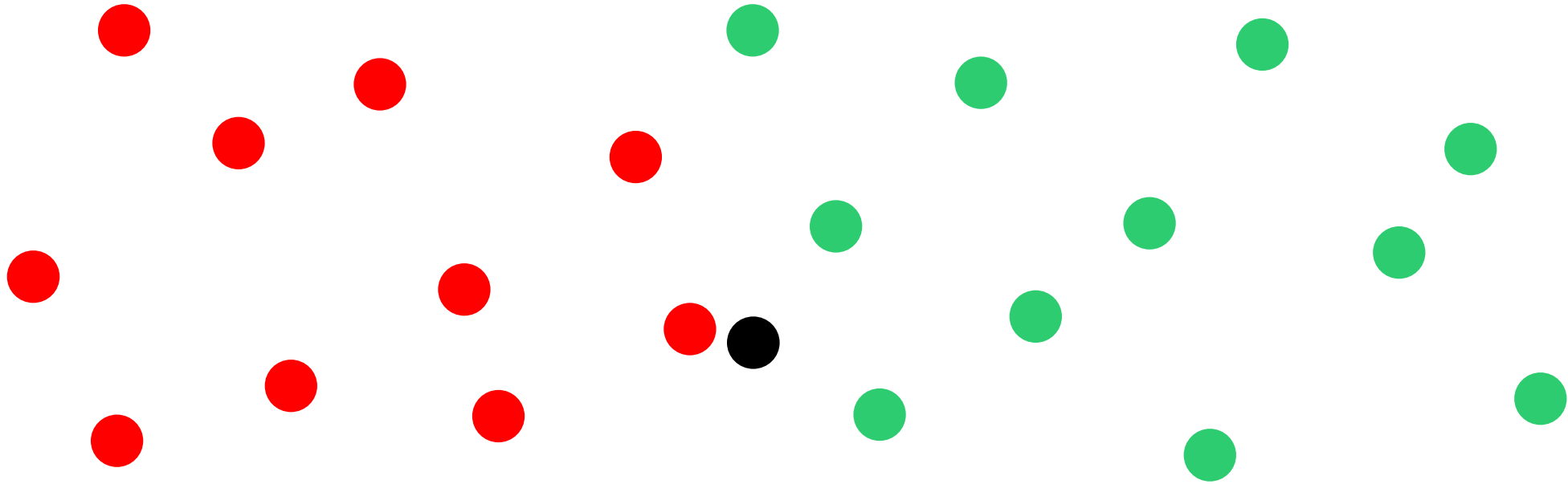
# Classification with Weighted rNN

26



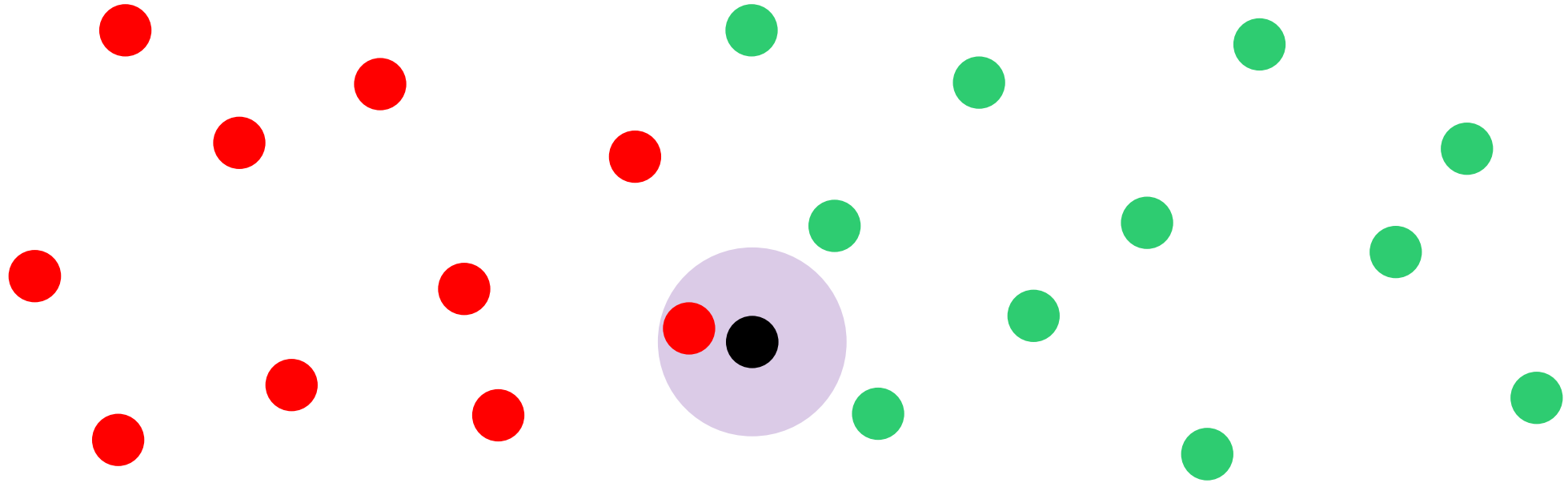
# Classification with Weighted rNN

26



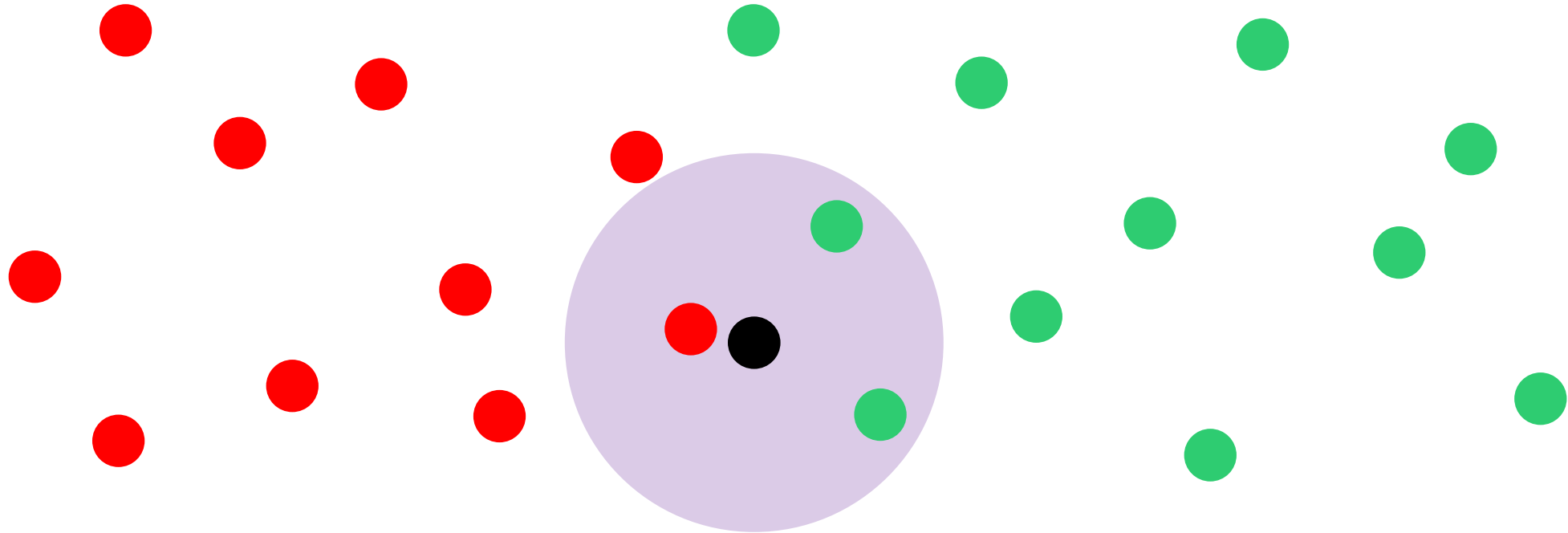
# Classification with Weighted rNN

26



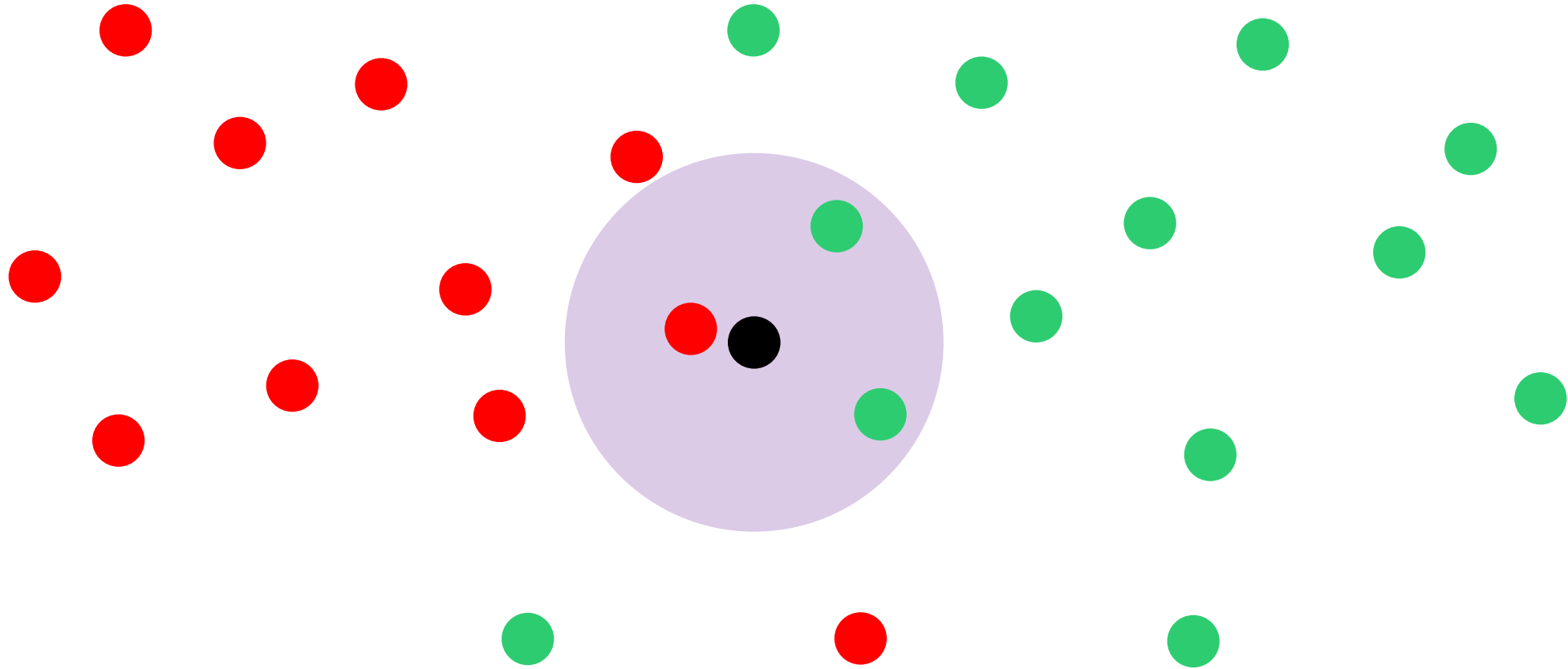
# Classification with Weighted rNN

26



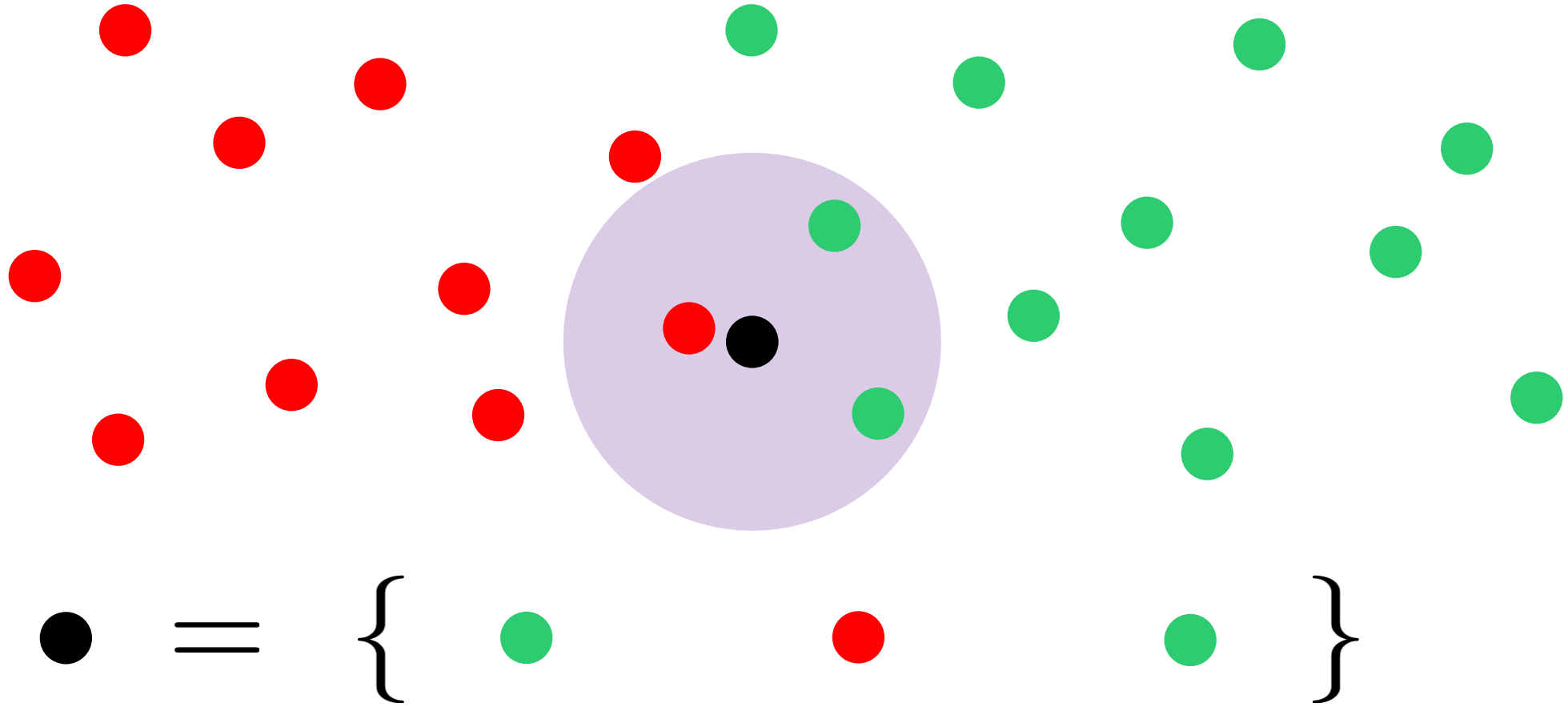
# Classification with Weighted rNN

26



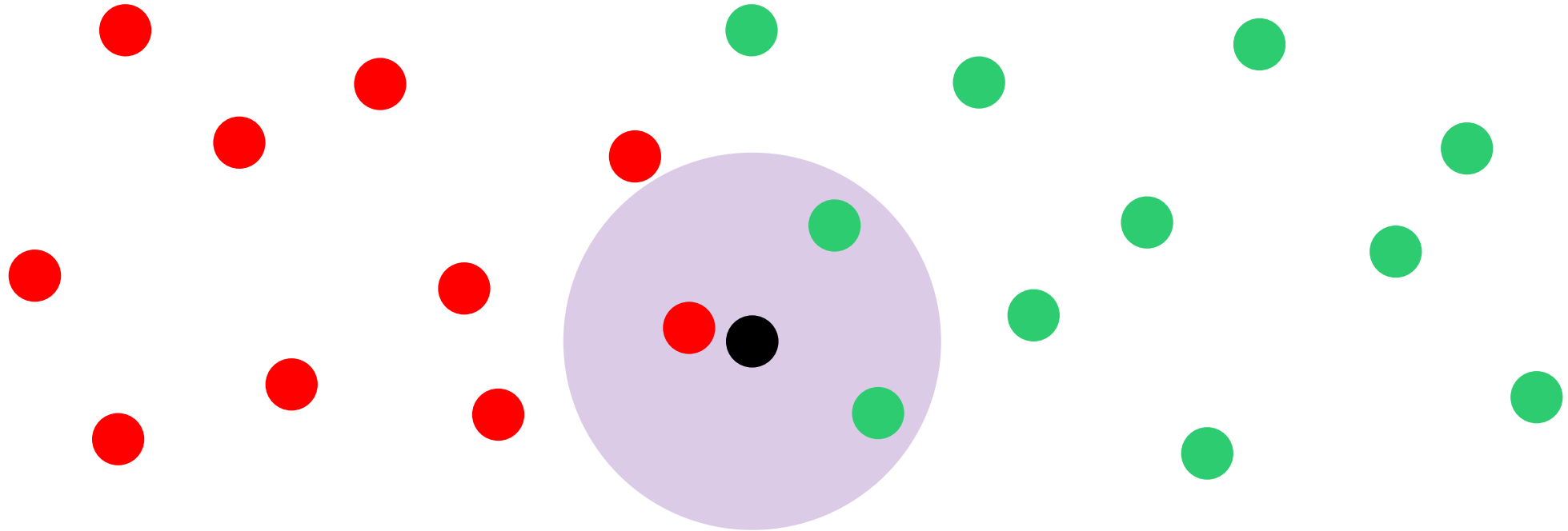
# Classification with Weighted rNN

26



# Classification with Weighted rNN

26

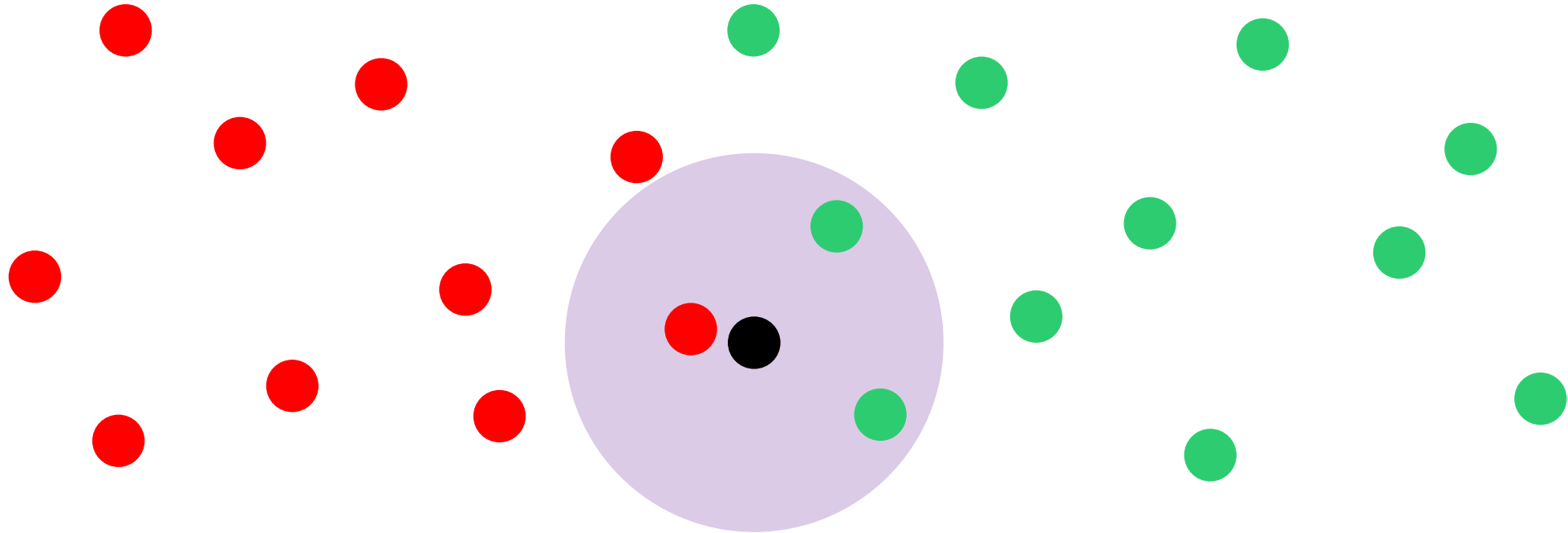


$$\bullet = \left\{ \frac{1}{3} \bullet \quad \frac{1}{3} \bullet \quad \frac{1}{3} \bullet \right\}$$

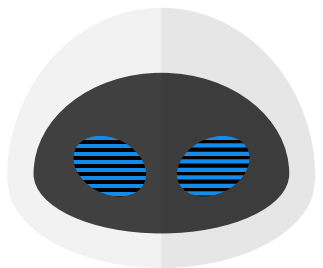


# Classification with Weighted rNN

26



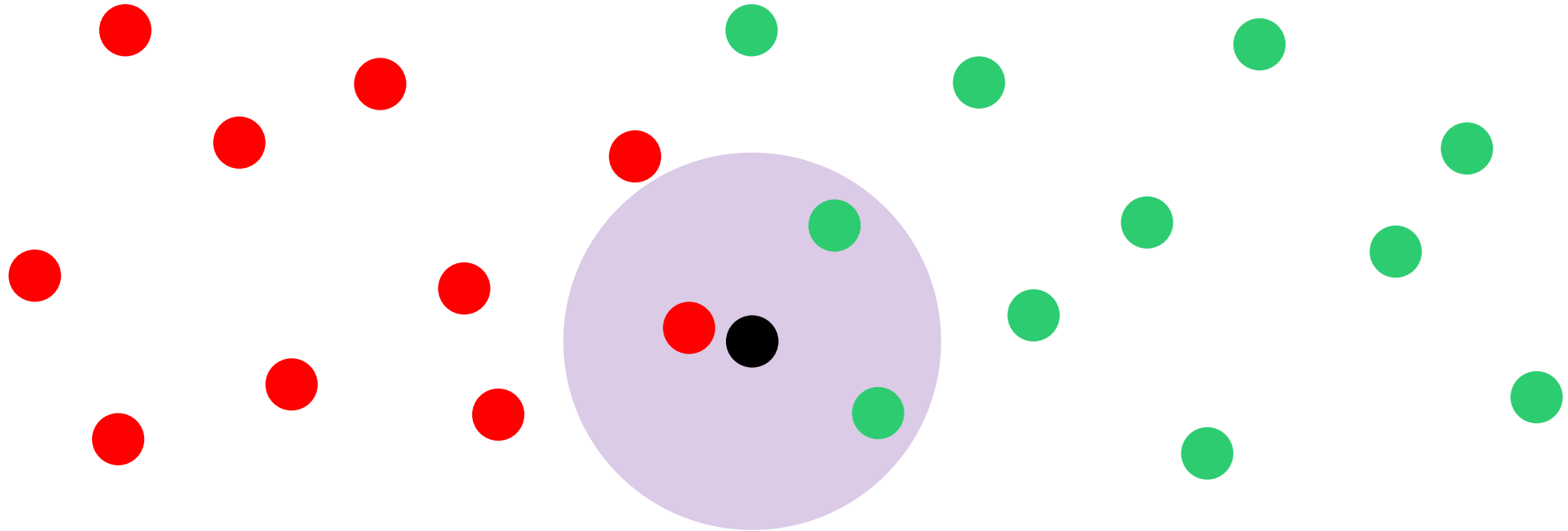
$$\bullet = \left\{ \frac{1}{3} \text{ (green dot)} \quad \frac{1}{3} \text{ (red dot)} \quad \frac{1}{3} \text{ (green dot)} \right\}$$





# Classification with Weighted rNN

26

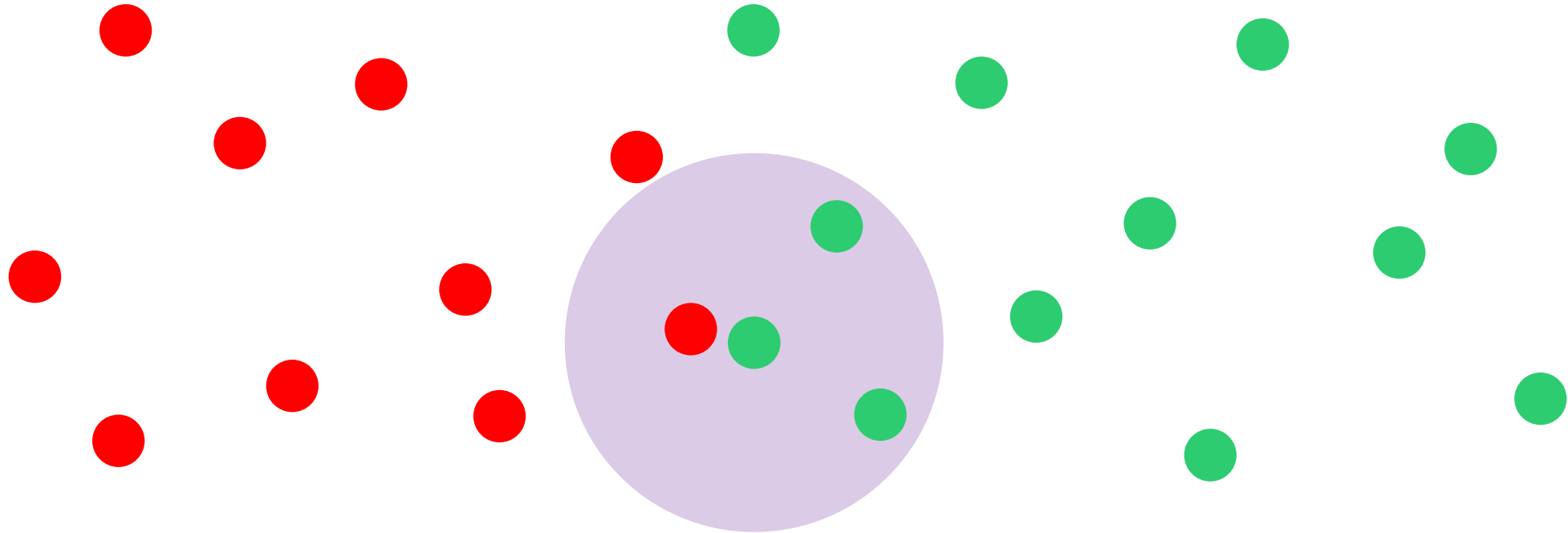


$$\bullet = \left\{ \frac{1}{3} \text{ green}, \frac{1}{3} \text{ red}, \frac{1}{3} \text{ green} \right\}$$

Green gets  $\frac{1}{3} + \frac{1}{3} = \frac{2}{3}$  votes whereas  
Red gets  $\frac{1}{3}$  votes – Green wins!!

# Classification with Weighted rNN

26

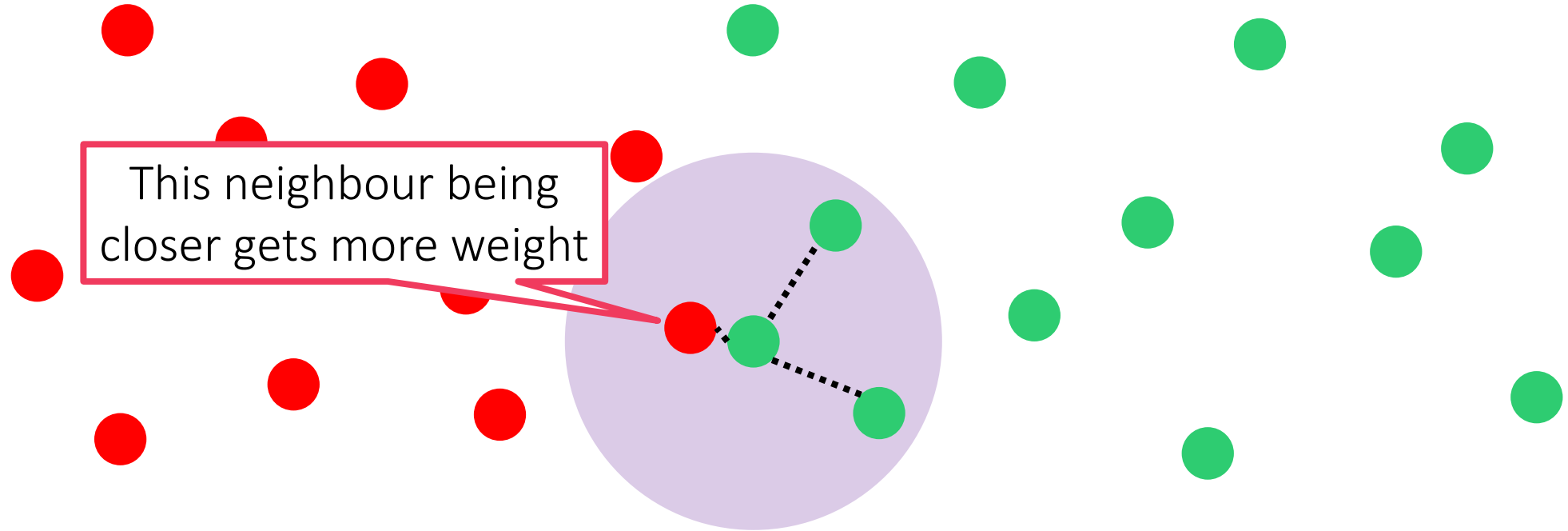


$$\text{Green} = \left\{ \frac{1}{3} \text{ Green} + \frac{1}{3} \text{ Red} + \frac{1}{3} \text{ Green} \right\}$$

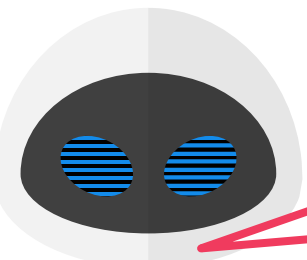
Green gets  $\frac{1}{3} + \frac{1}{3} = \frac{2}{3}$  votes whereas  
Red gets  $\frac{1}{3}$  votes – Green wins!!

# Classification with Weighted rNN

26



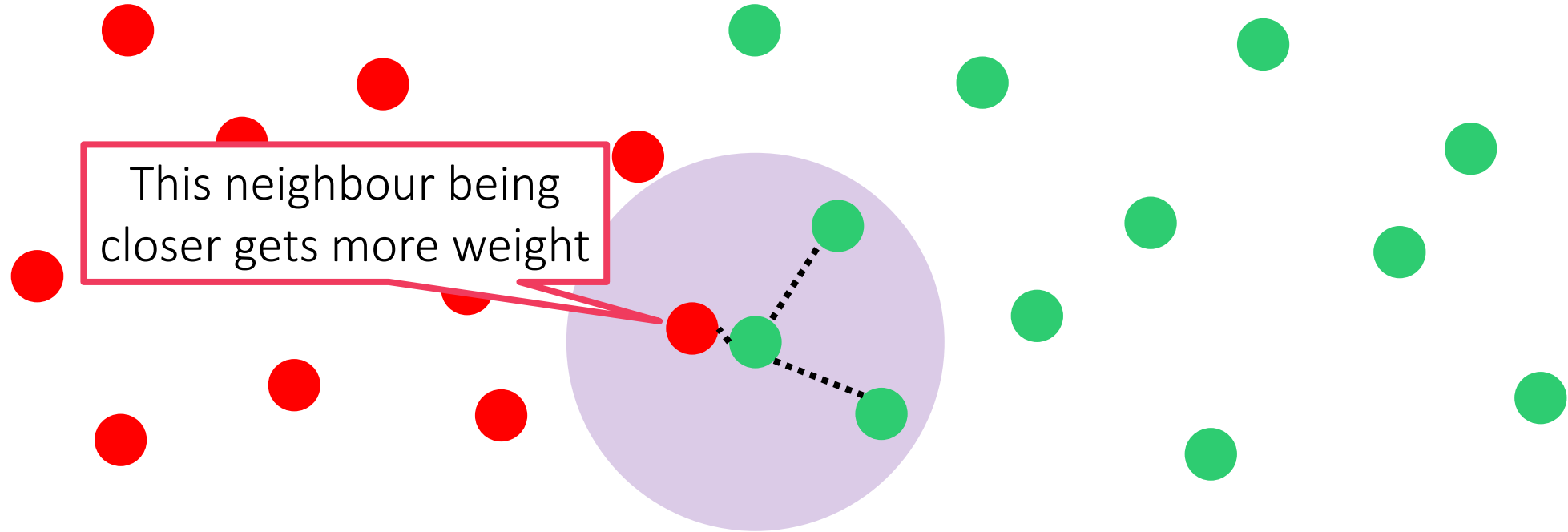
$$\text{Green} = \left\{ \frac{1}{3} \text{ Green} + \frac{1}{3} \text{ Red} + \frac{1}{3} \text{ Green} \right\}$$



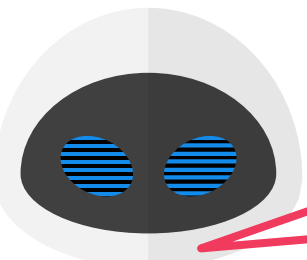
Green gets  $\frac{1}{3} + \frac{1}{3} = \frac{2}{3}$  votes whereas Red gets  $\frac{1}{3}$  votes – Green wins!!

# Classification with Weighted rNN

26



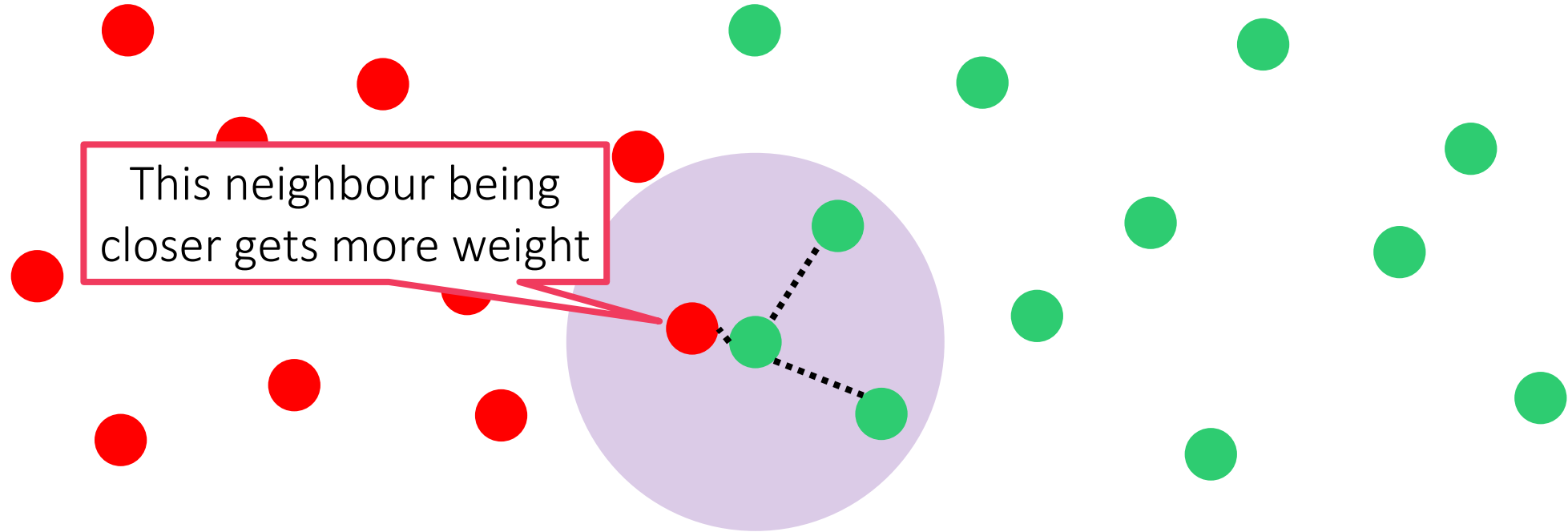
$$\text{Green} = \left\{ \frac{1}{5} \text{ Green} + \frac{3}{5} \text{ Red} + \frac{1}{5} \text{ Green} \right\}$$



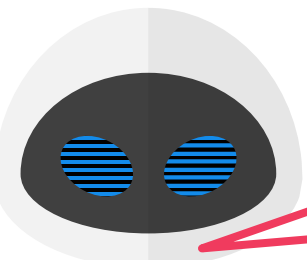
Green gets  $\frac{1}{3} + \frac{1}{3} = \frac{2}{3}$  votes whereas Red gets  $\frac{1}{3}$  votes – Green wins!!

# Classification with Weighted rNN

26



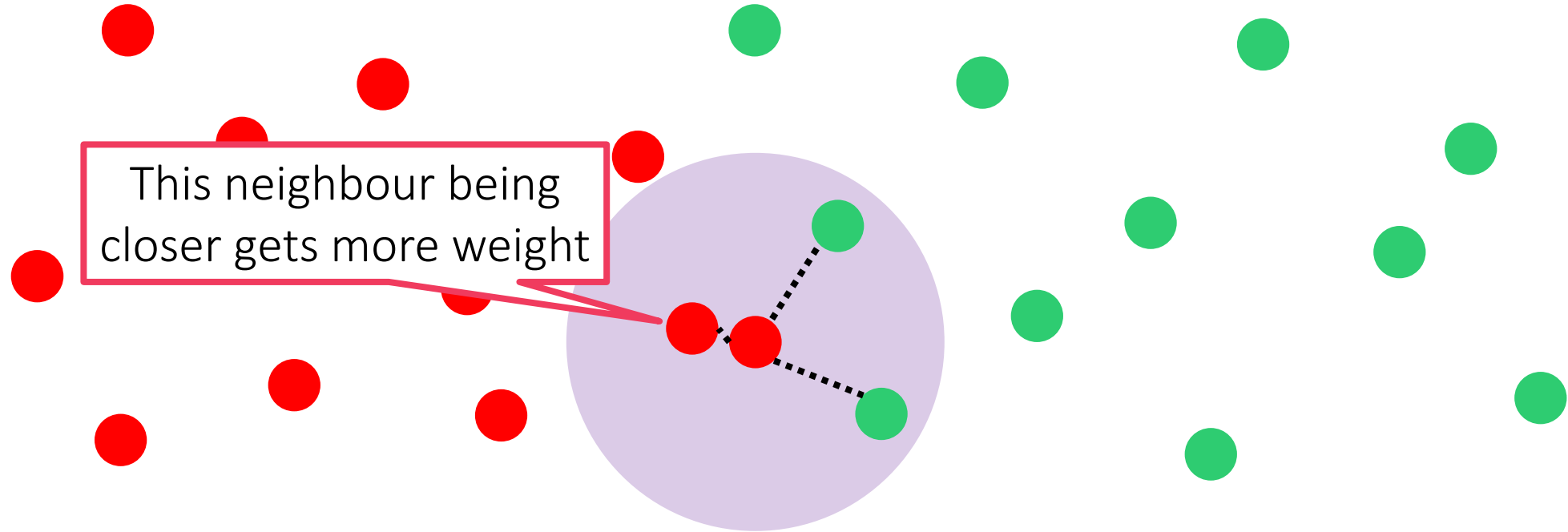
$$\text{Green} = \left\{ \frac{1}{5} \text{ Green} + \frac{3}{5} \text{ Red} + \frac{1}{5} \text{ Green} \right\}$$



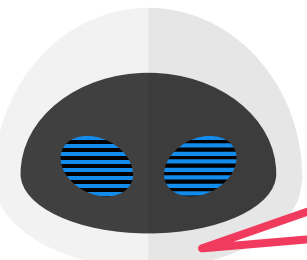
Green gets  $\frac{1}{5} + \frac{1}{5} = \frac{2}{5}$  votes whereas Red gets  $\frac{3}{5}$  votes – Red wins!!

# Classification with Weighted rNN

26



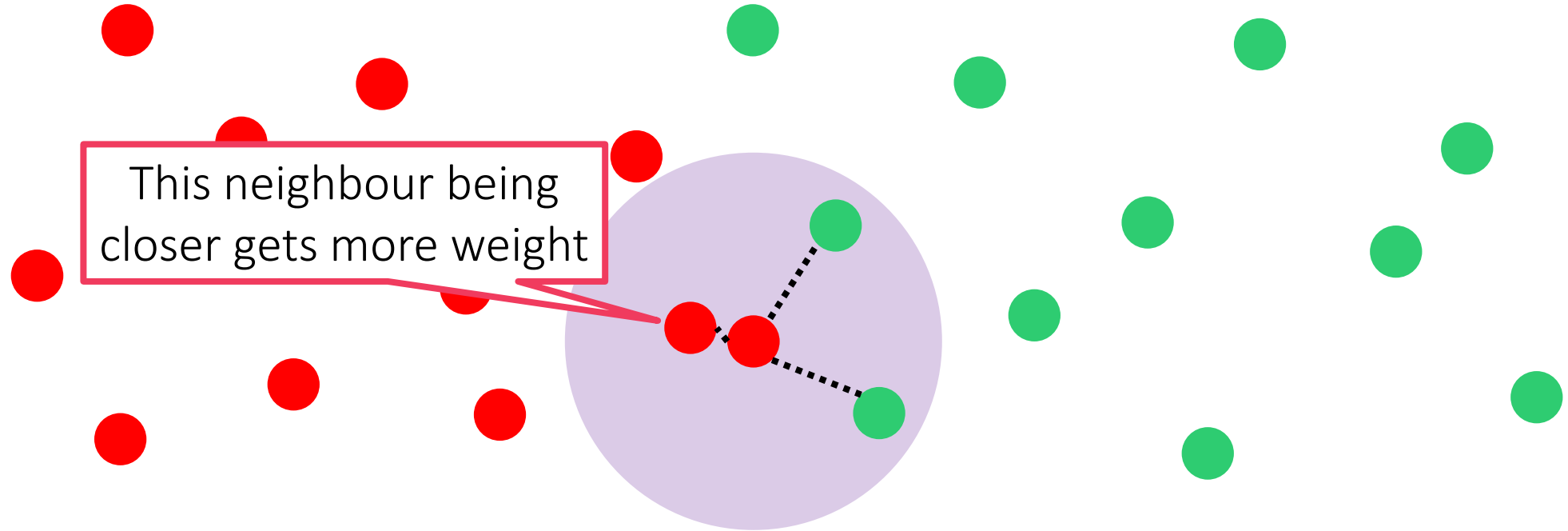
$$\text{Red} = \left\{ \frac{1}{5} \text{ Green} \quad \frac{3}{5} \text{ Red} \quad \frac{1}{5} \text{ Green} \right\}$$



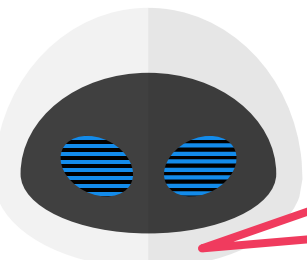
Green gets  $\frac{1}{5} + \frac{1}{5} = \frac{2}{5}$  votes whereas  
Red gets  $\frac{3}{5}$  votes – Red wins!!

# Classification with Weighted rNN

26

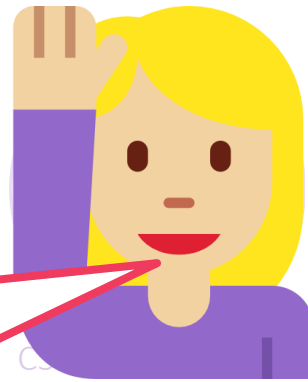


$$\text{Red} = \left\{ \frac{1}{5} \text{ Green} \quad \frac{3}{5} \text{ Red} \quad \frac{1}{5} \text{ Green} \right\}$$



Green gets  $\frac{1}{5} + \frac{1}{5} = \frac{2}{5}$  votes whereas Red gets  $\frac{3}{5}$  votes – Red wins!!

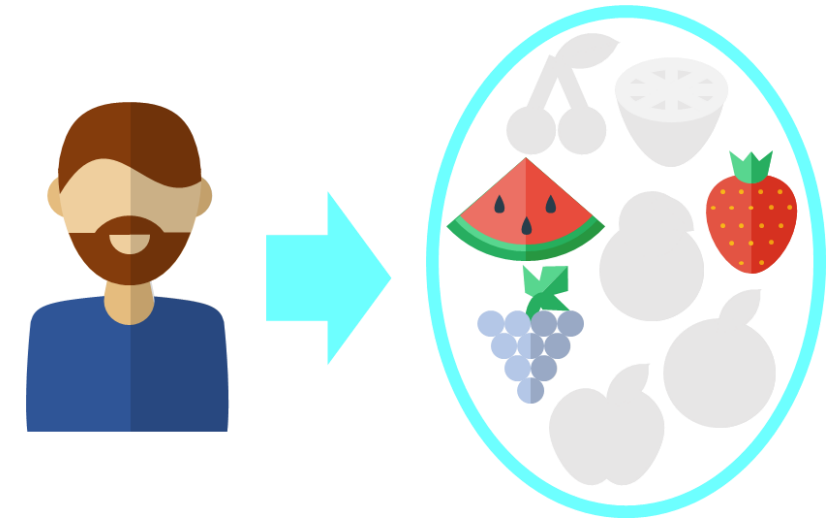
This method elegantly works even if there are more than 2 classes!



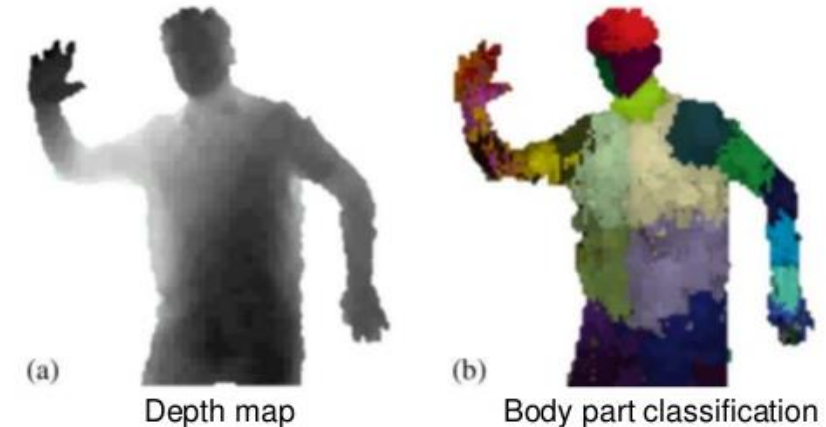
# Decision Trees

- Give a *faster* way to perform NN search
- Can handle non-numeric features too!
- Very popular in ML – classification, recsys
- Extremely fast at making predictions
- Easy to interpret by humans too
- Model size can be very large ☹️
- Can give good train perf. but bad test perf. (known as *overfitting*)

Recommendation Systems



Gaming – Kinect for Xbox 360





# Decision Trees for Classification

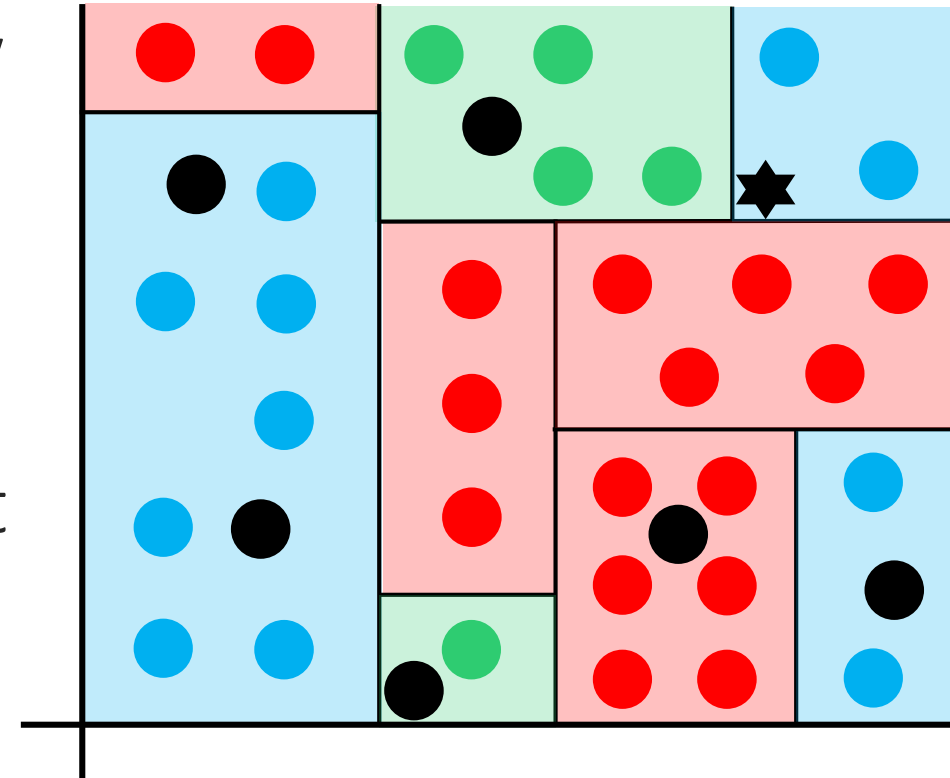
41

Finding the nearest neighbours can be slow

So repeatedly partition feature space  $\mathbb{R}^d$

For test data point, much faster to find out which partition is it in which it lies

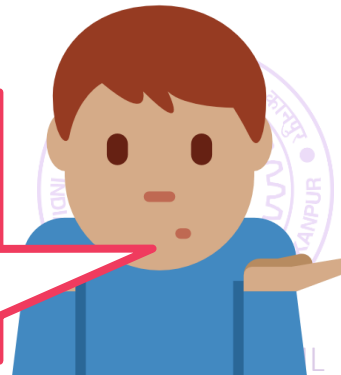
Can consider the other (train) points in that partition as approximate neighbours



However, if your partitions are fine enough, this will happen very rarely and not hurt performance too much!

Yes, decision trees may not always give you the exact neighbors for points lying at boundary of a partition

What about points like this?  
The nearest neighbours are lying in another partition ☹️



# Decision Trees for Classification

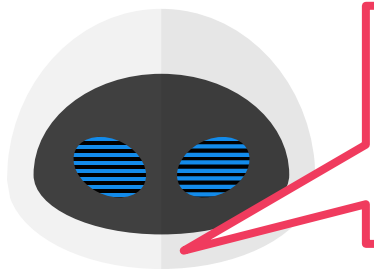
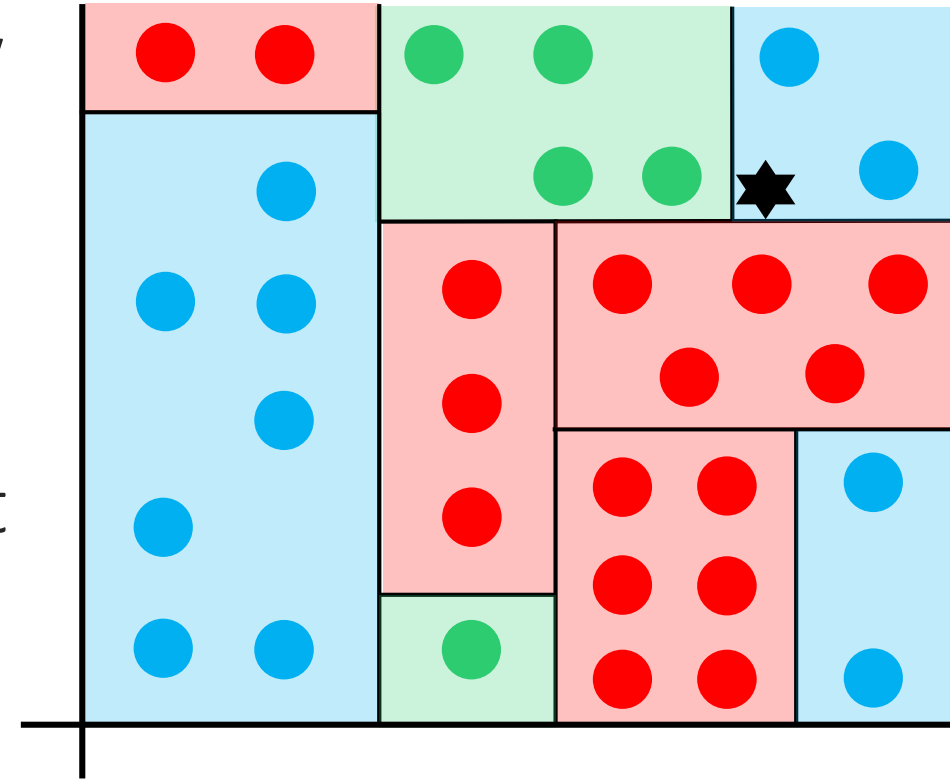
42

Finding the nearest neighbours can be slow

So repeatedly partition feature space  $\mathbb{R}^d$

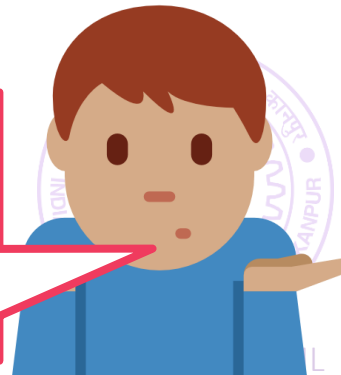
For test data point, much faster to find out which partition is it in which it lies

Can consider the other (train) points in that partition as approximate neighbours



Yes, decision trees may not always give you the exact neighbors for points lying at boundary of a partition

What about points like this?  
The nearest neighbours are lying in another partition ☹️



# Decision Trees for Classification

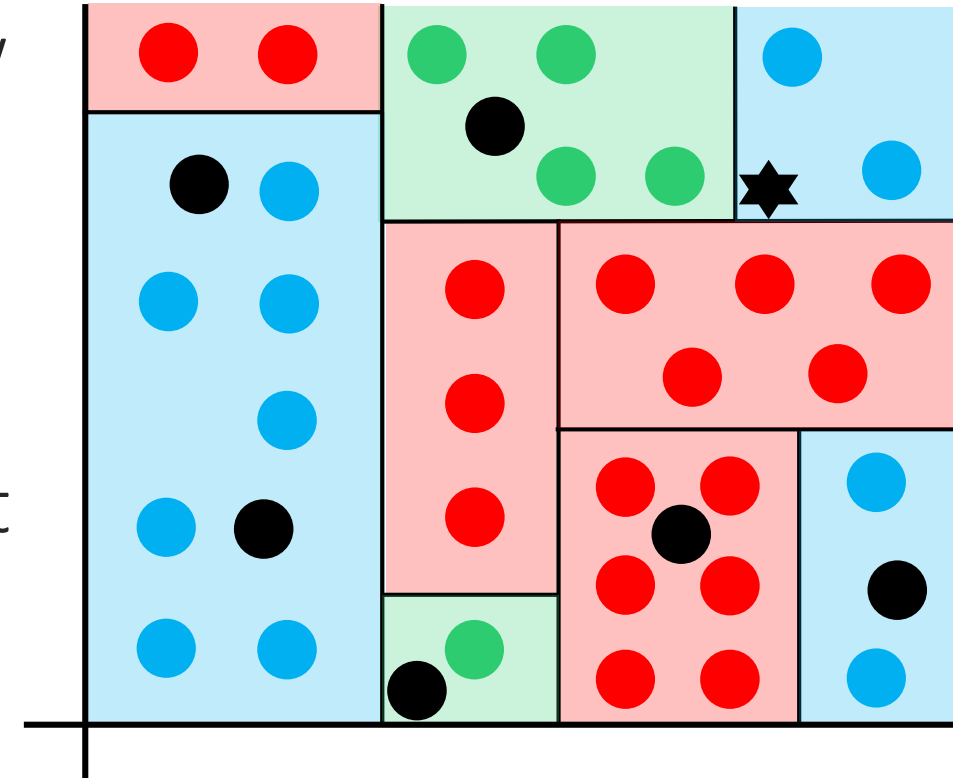
43

Finding the nearest neighbours can be slow

So repeatedly partition feature space  $\mathbb{R}^d$

For test data point, much faster to find out which partition is it in which it lies

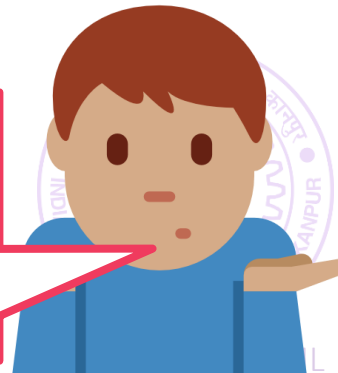
Can consider the other (train) points in that partition as approximate neighbours



However, if your partitions are fine enough, this will happen very rarely and not hurt performance too much!

Yes, decision trees may not always give you the exact neighbors for points lying at boundary of a partition

What about points like this?  
The nearest neighbours are lying in another partition ☹️



# Building Decision Trees

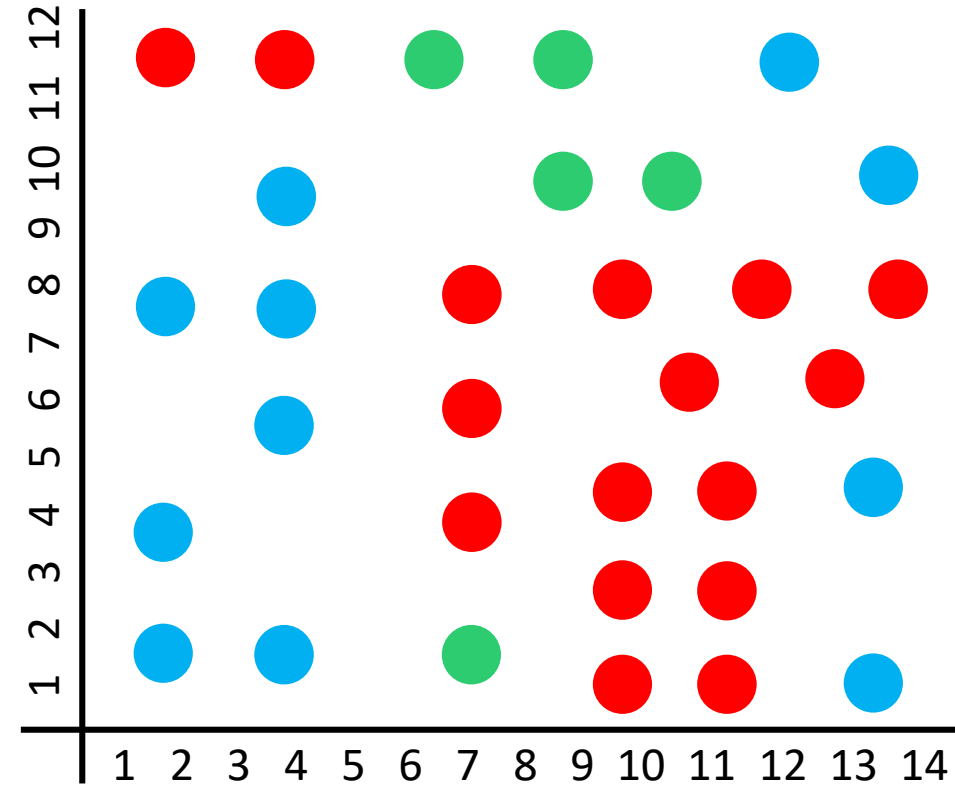
44



# Building Decision Trees

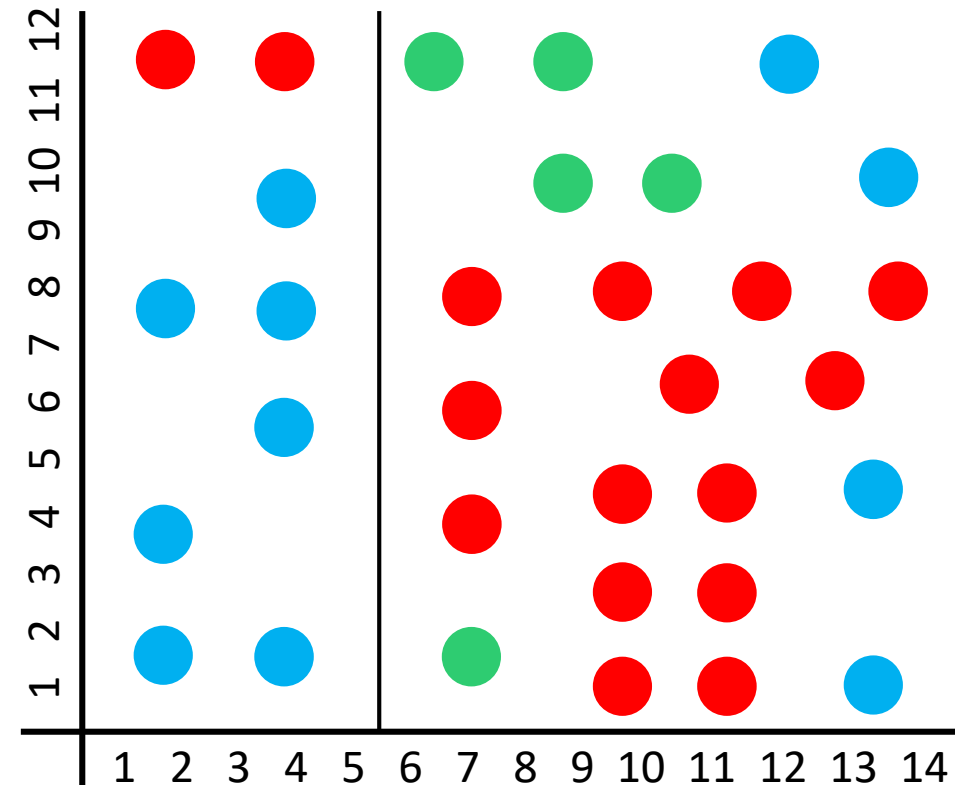
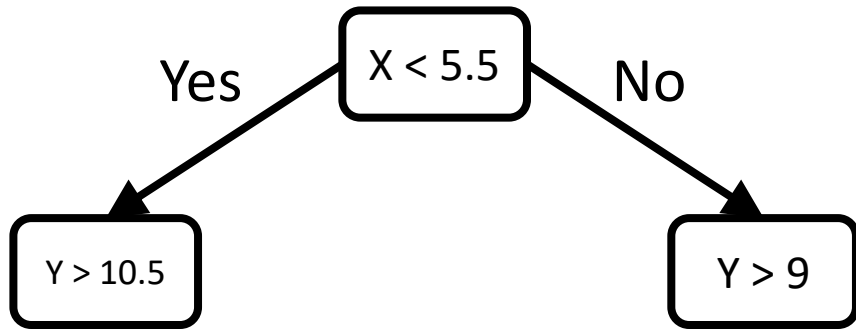
44

$X < 5.5$



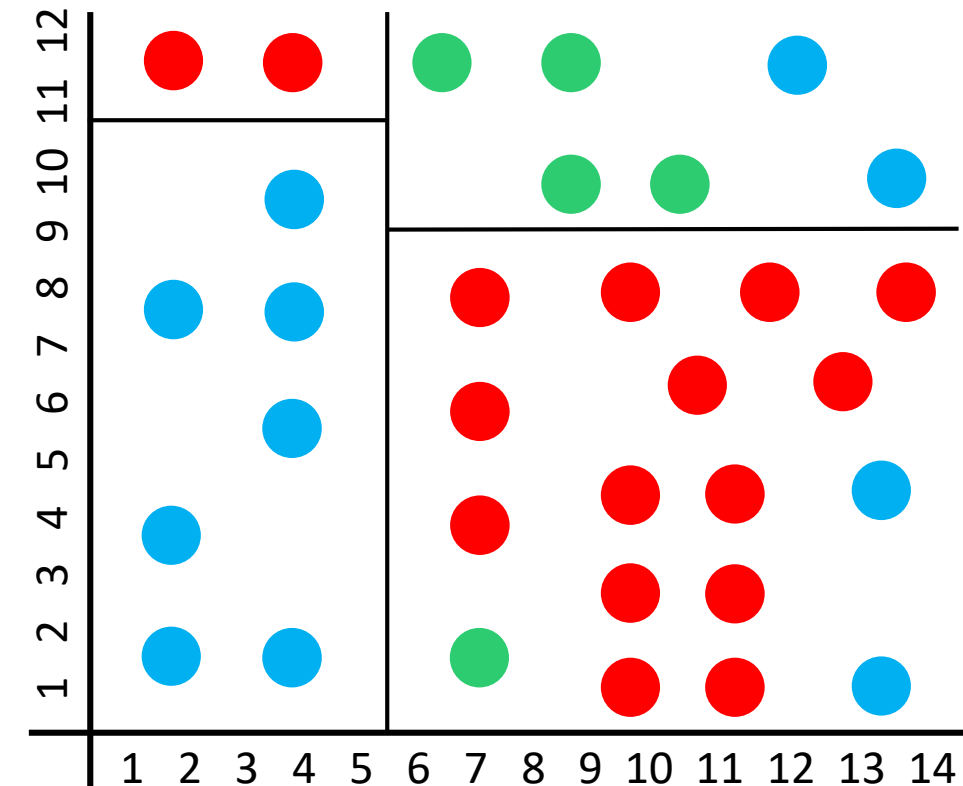
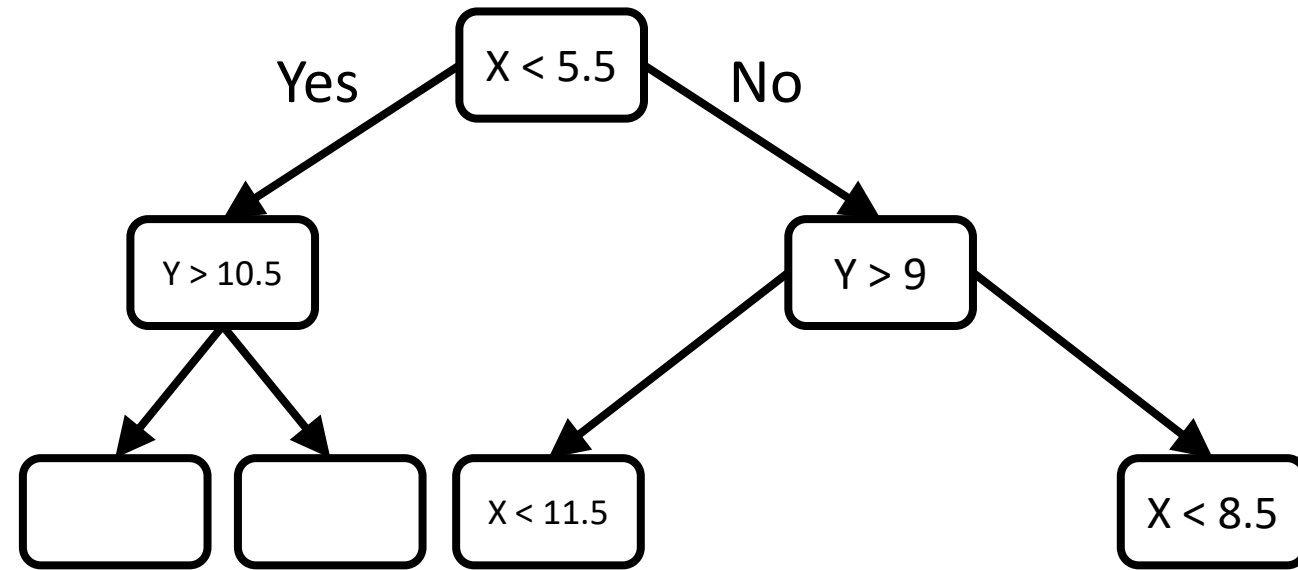
# Building Decision Trees

44



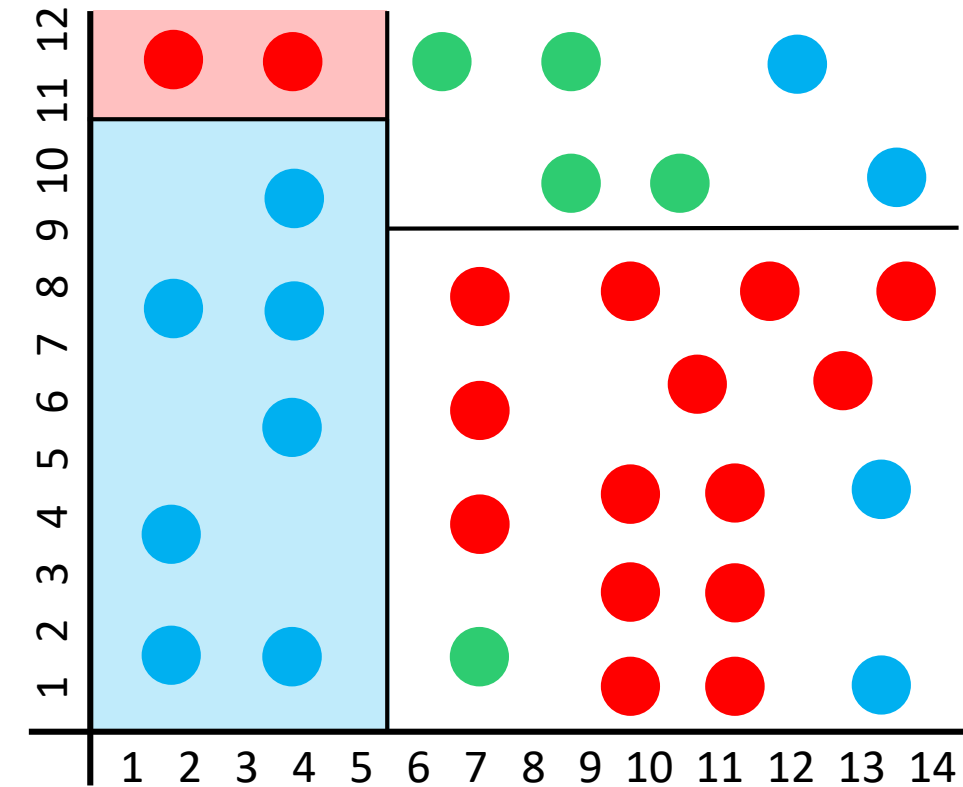
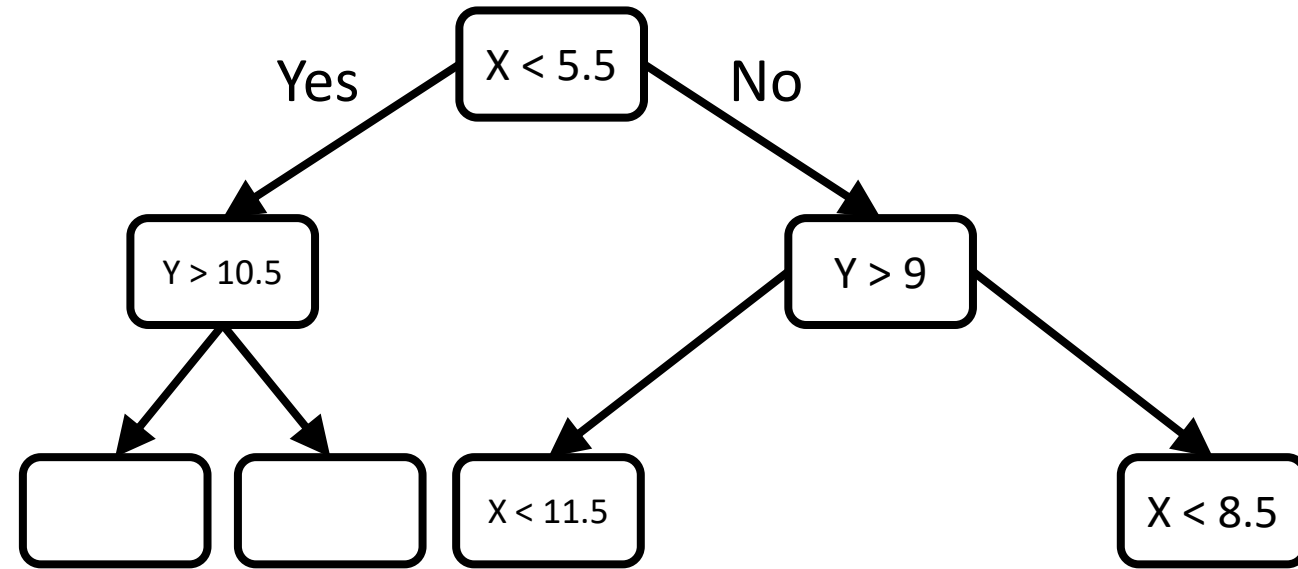
# Building Decision Trees

44



# Building Decision Trees

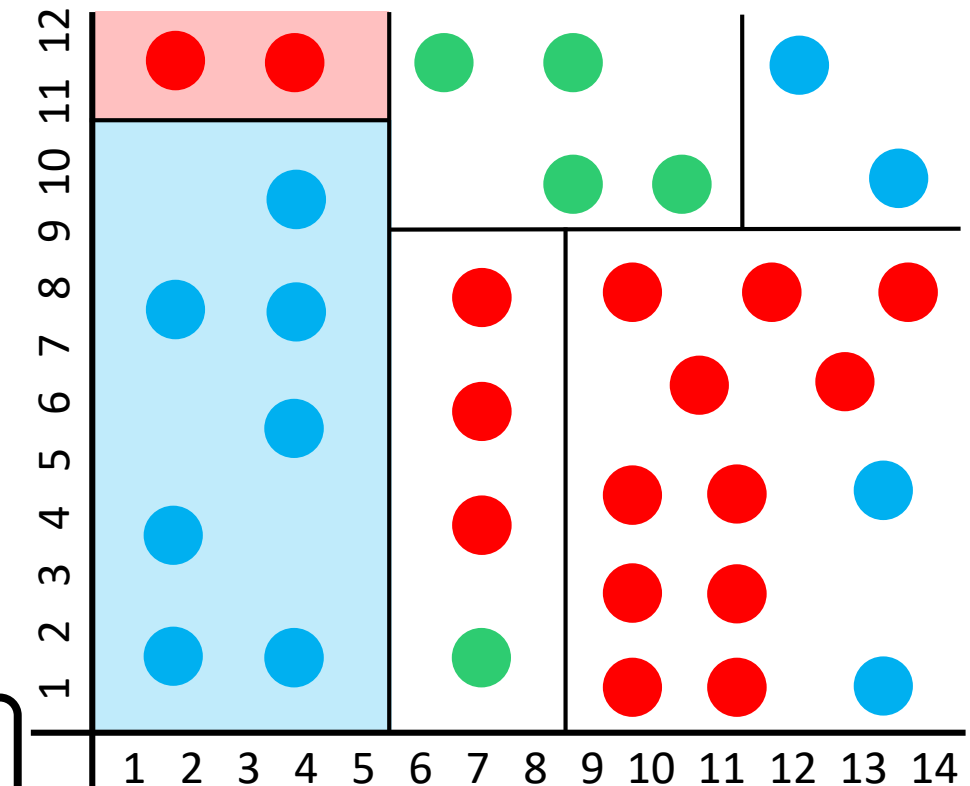
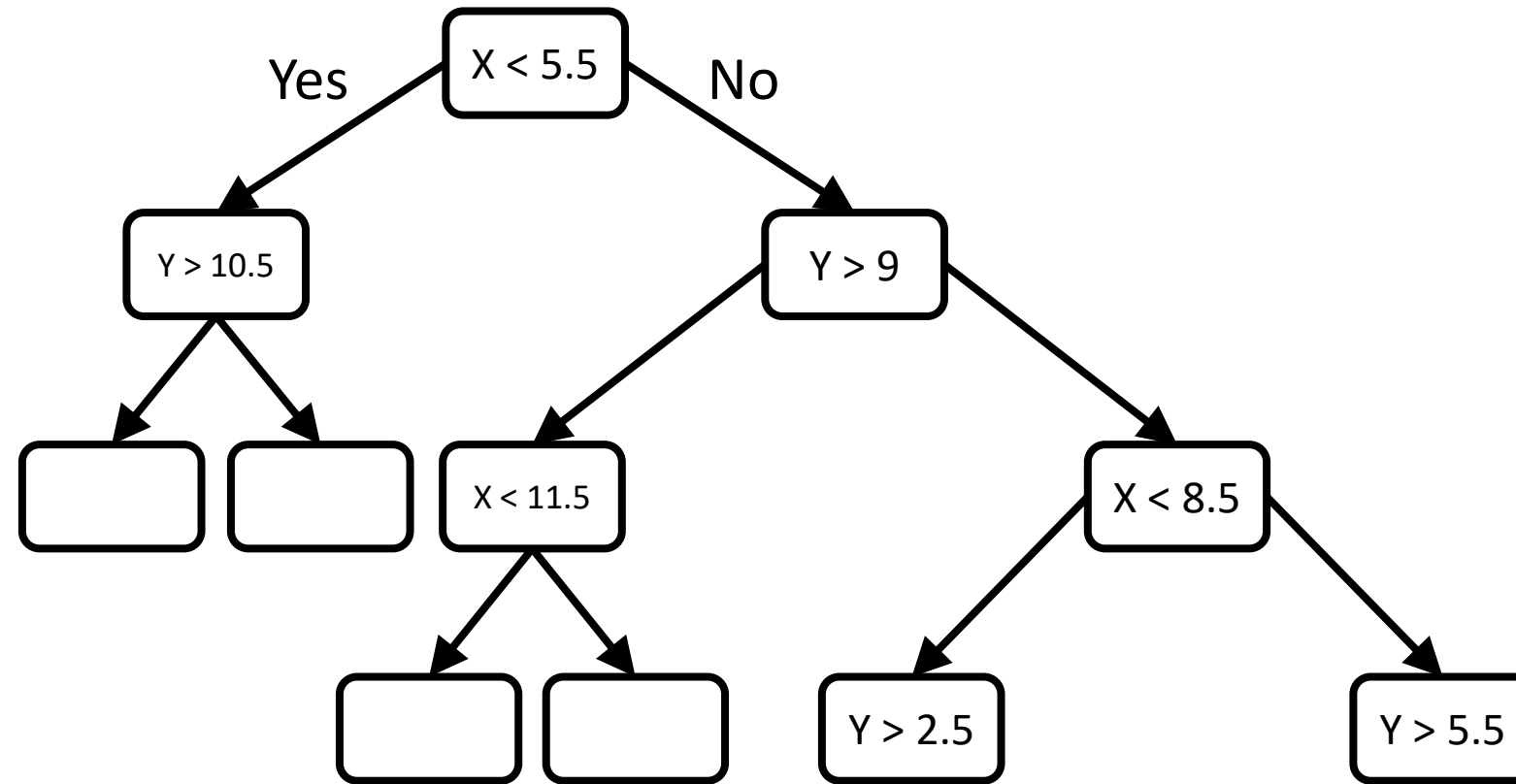
44





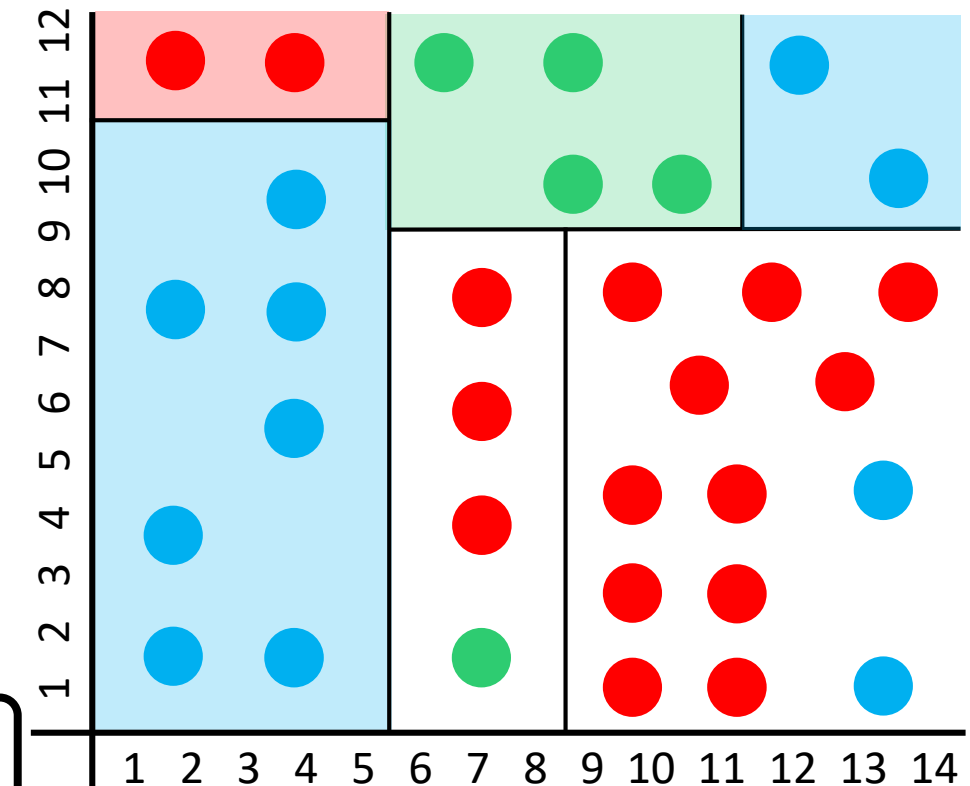
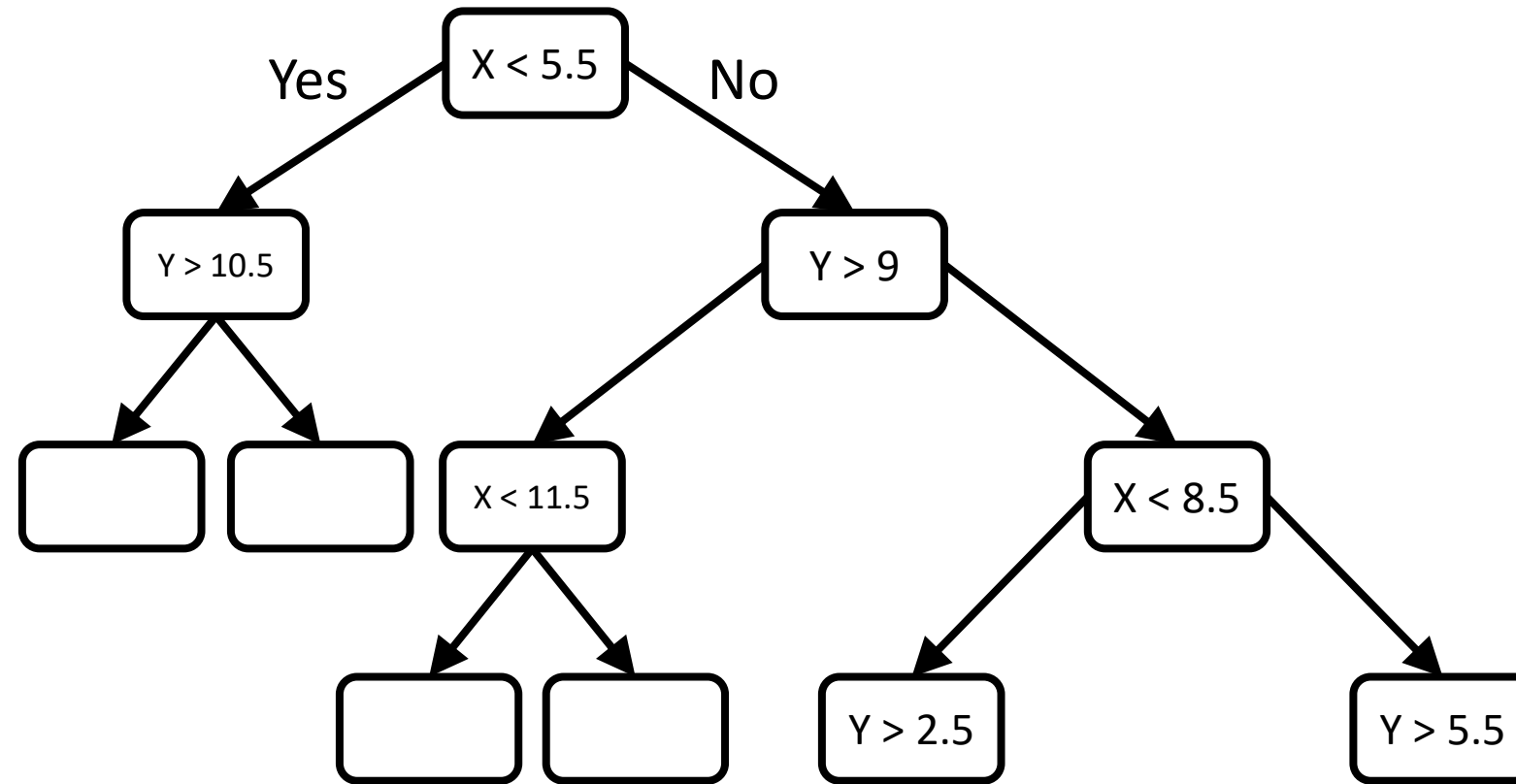
# Building Decision Trees

44



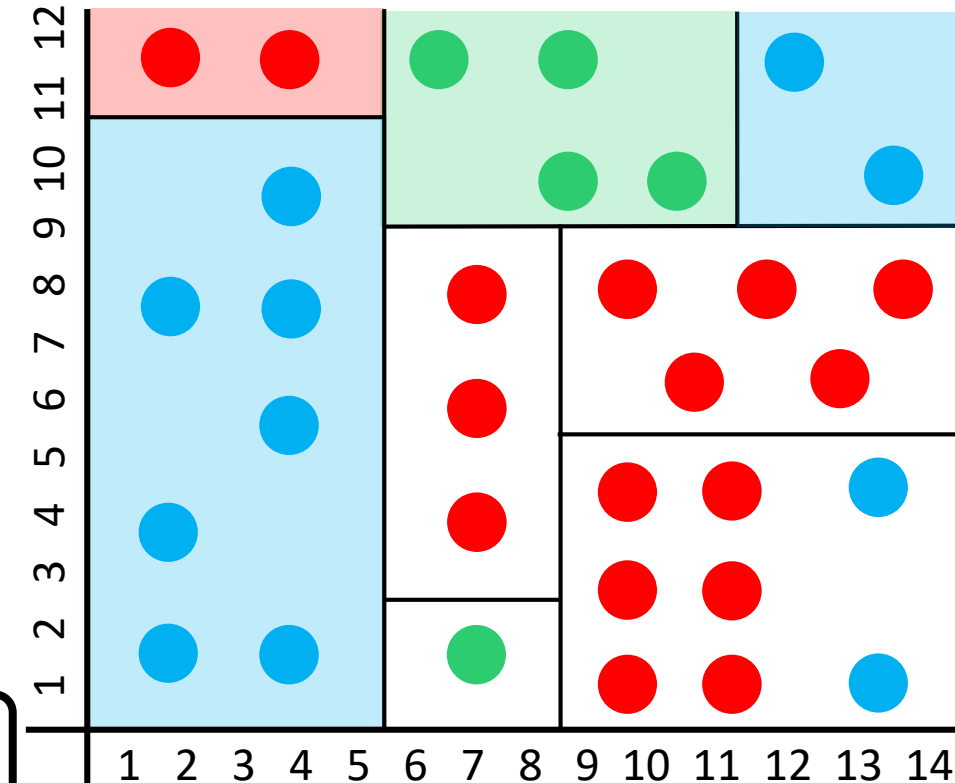
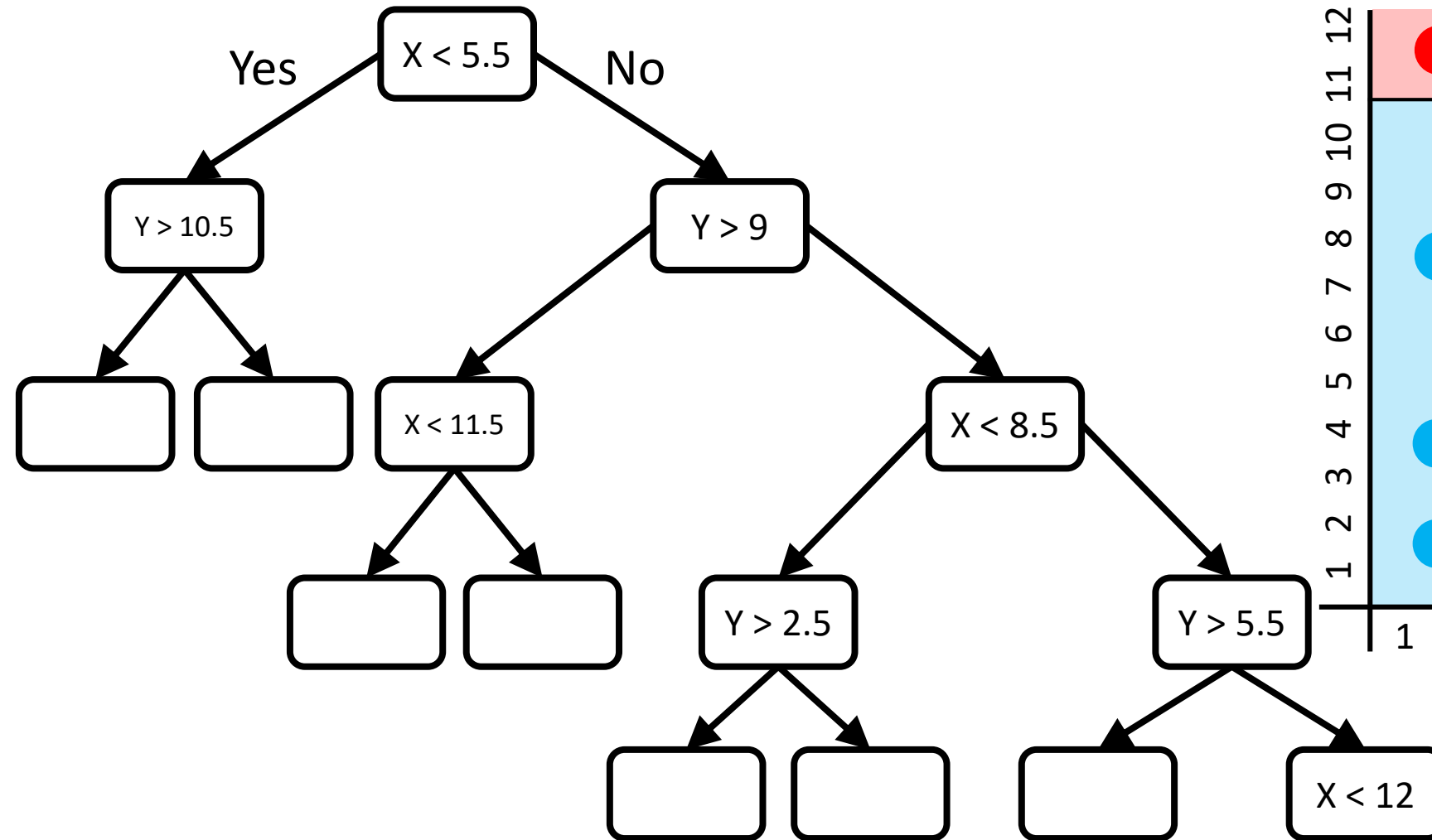
# Building Decision Trees

44



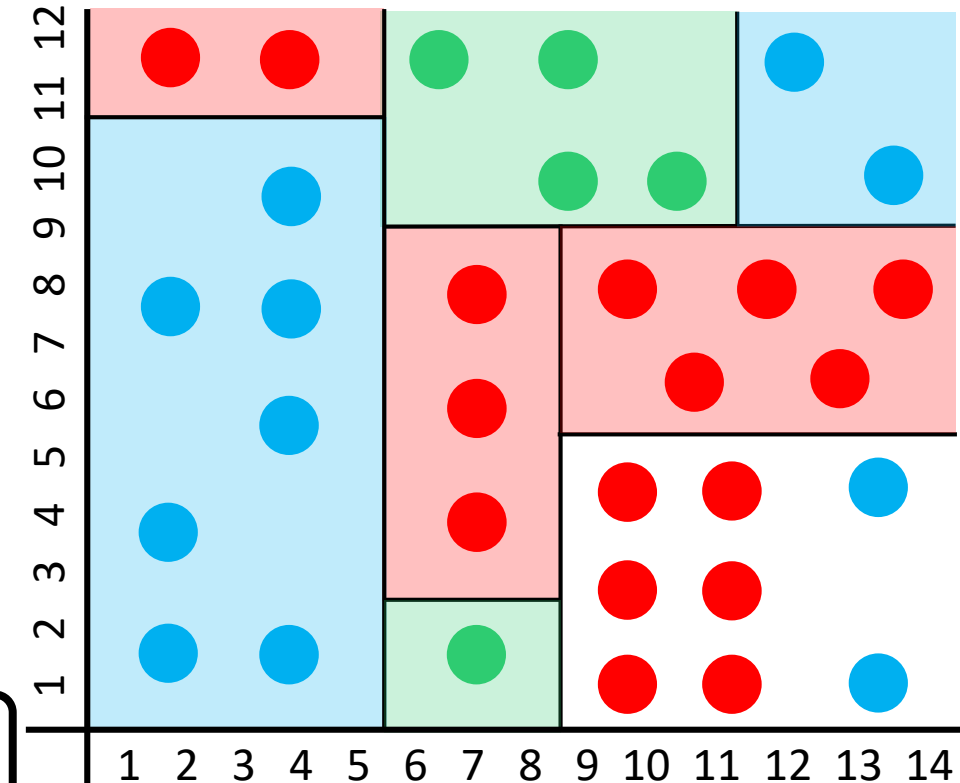
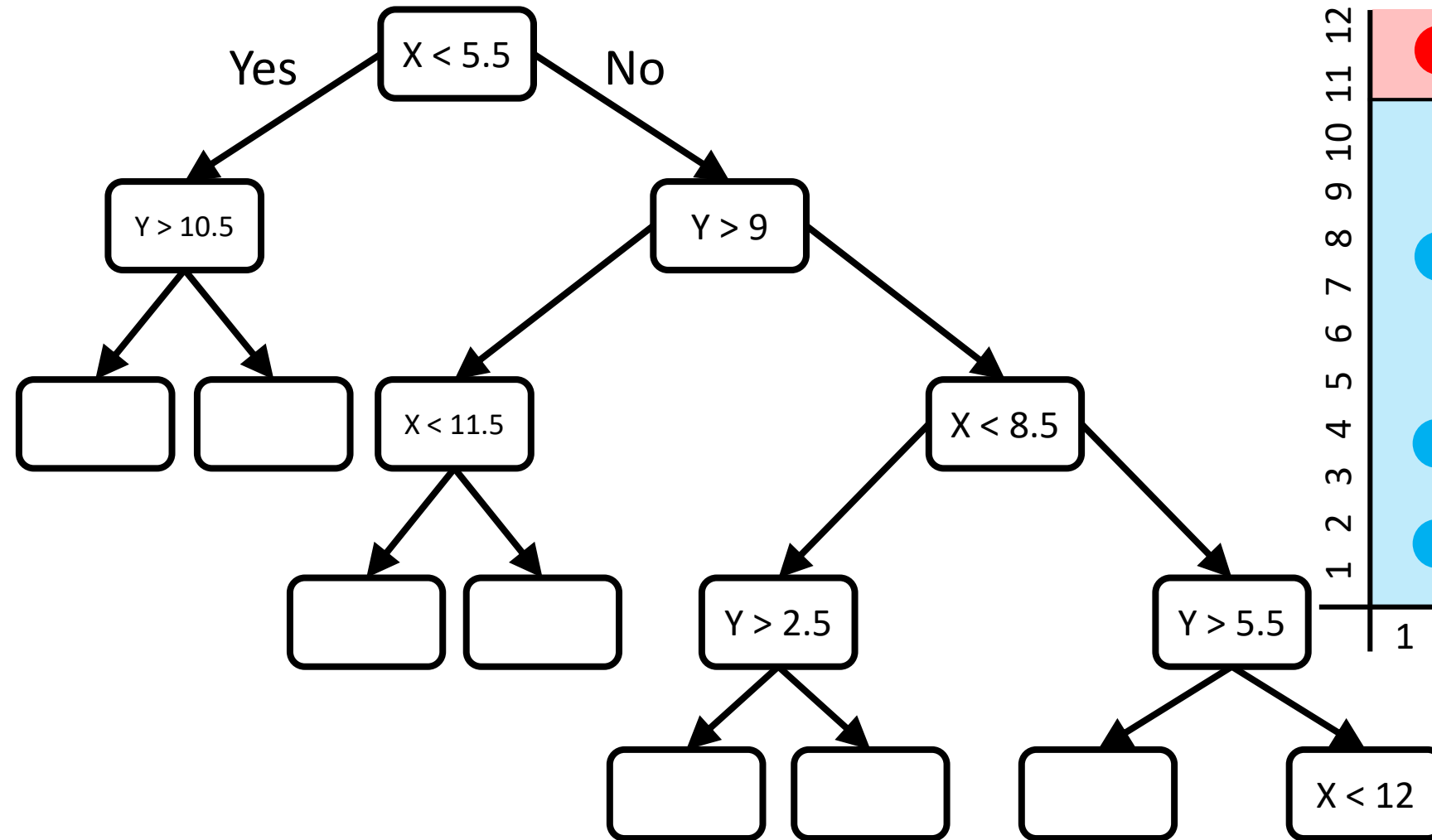
# Building Decision Trees

44



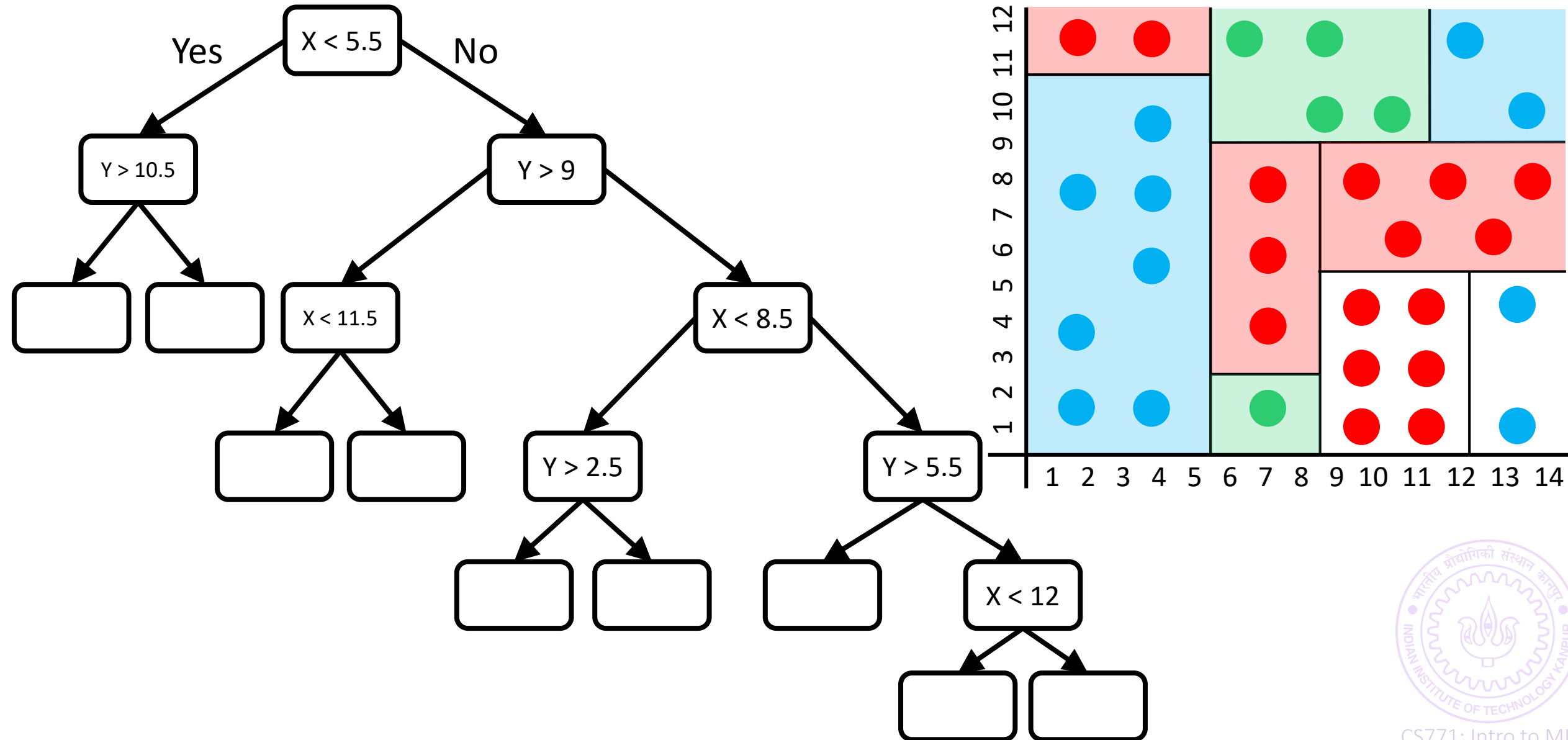
# Building Decision Trees

44



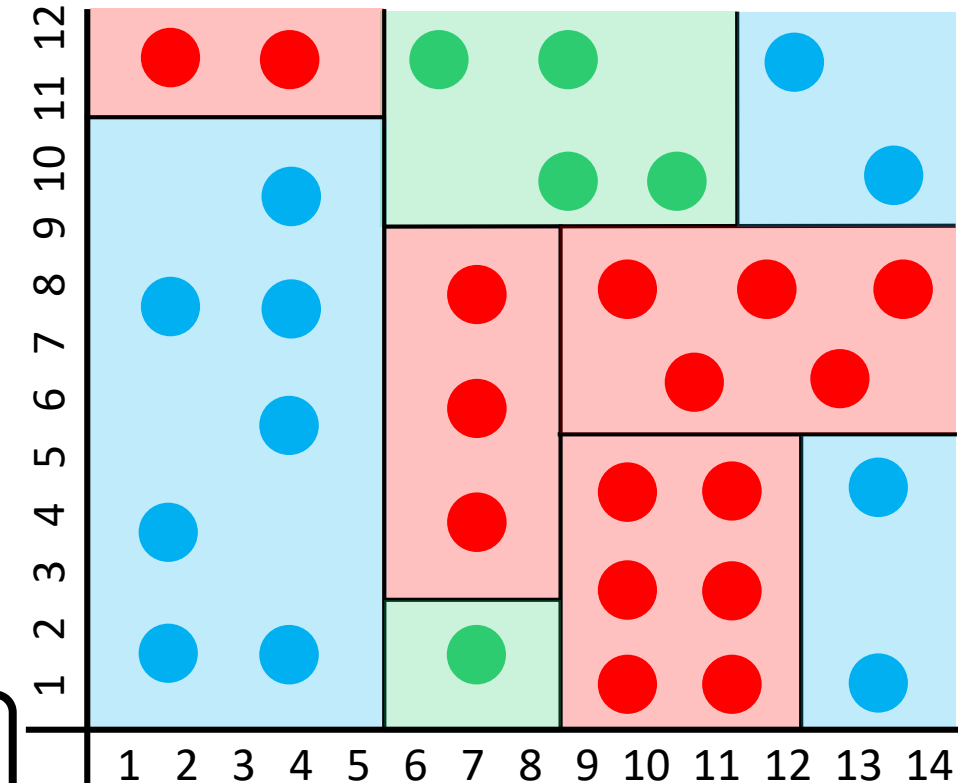
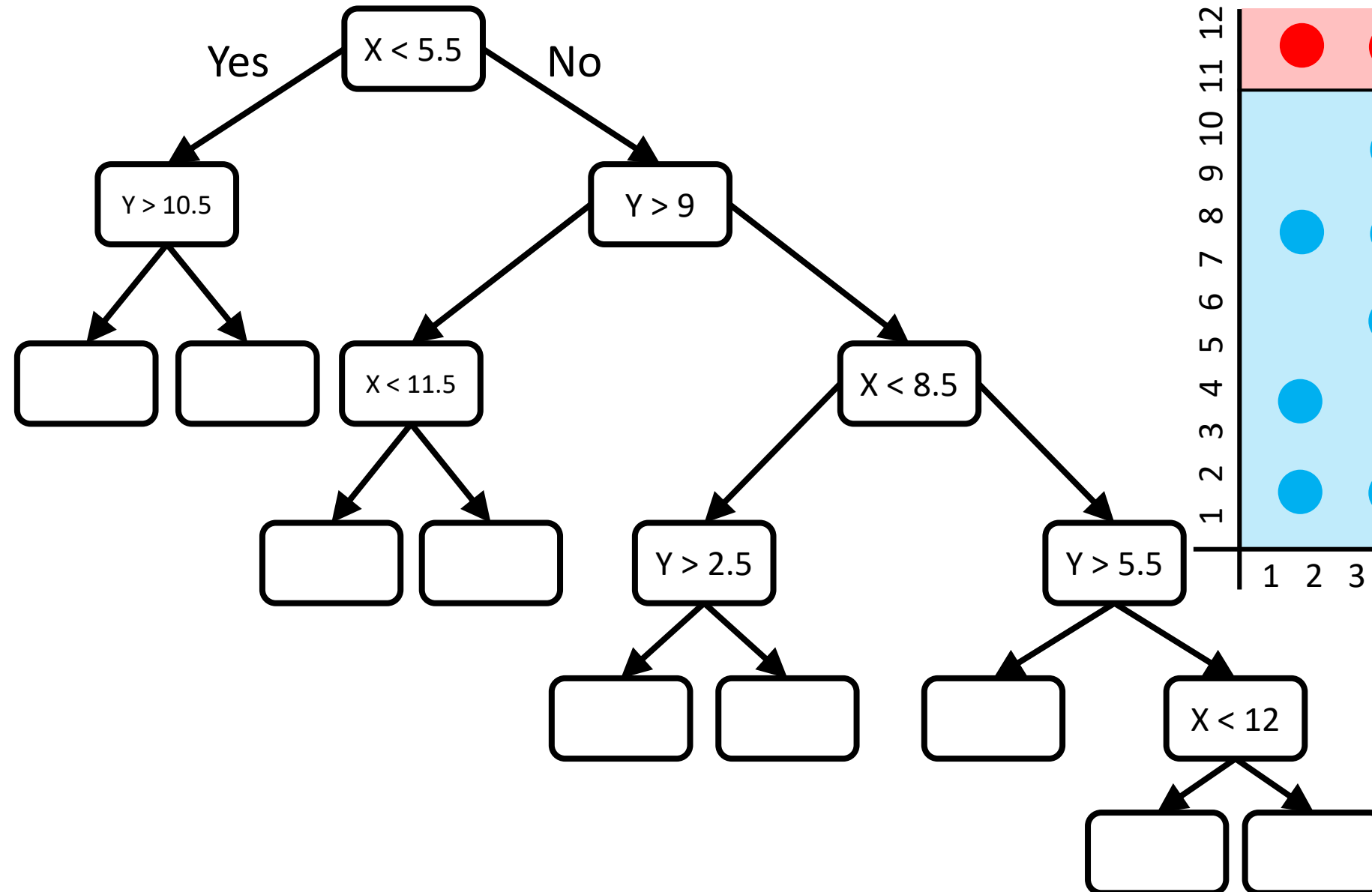
# Building Decision Trees

44



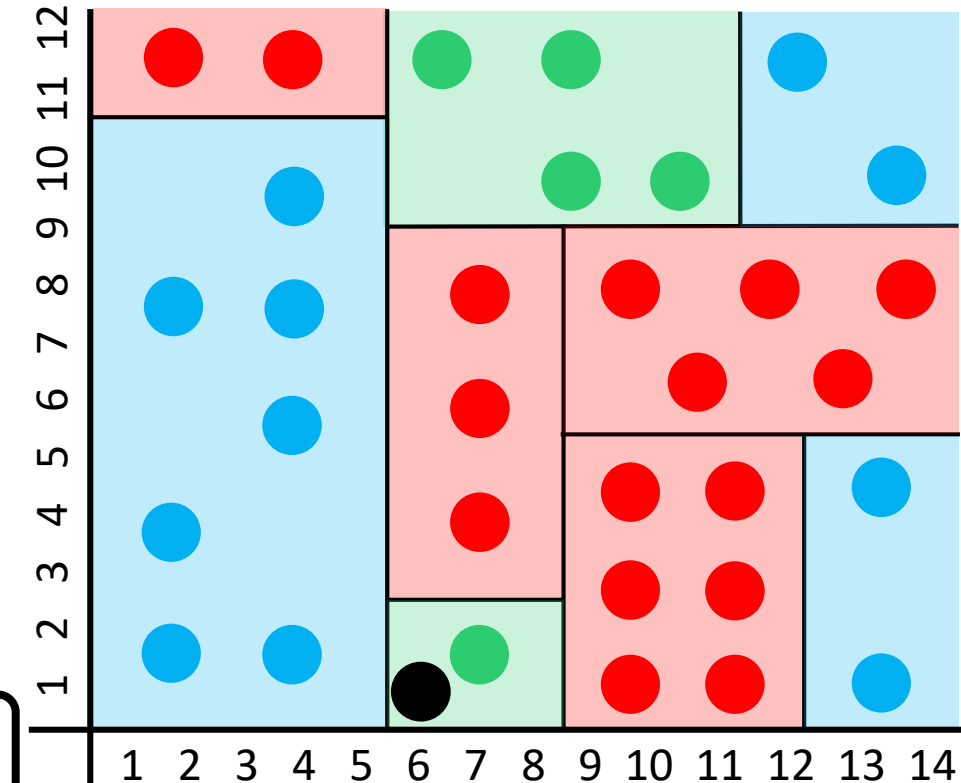
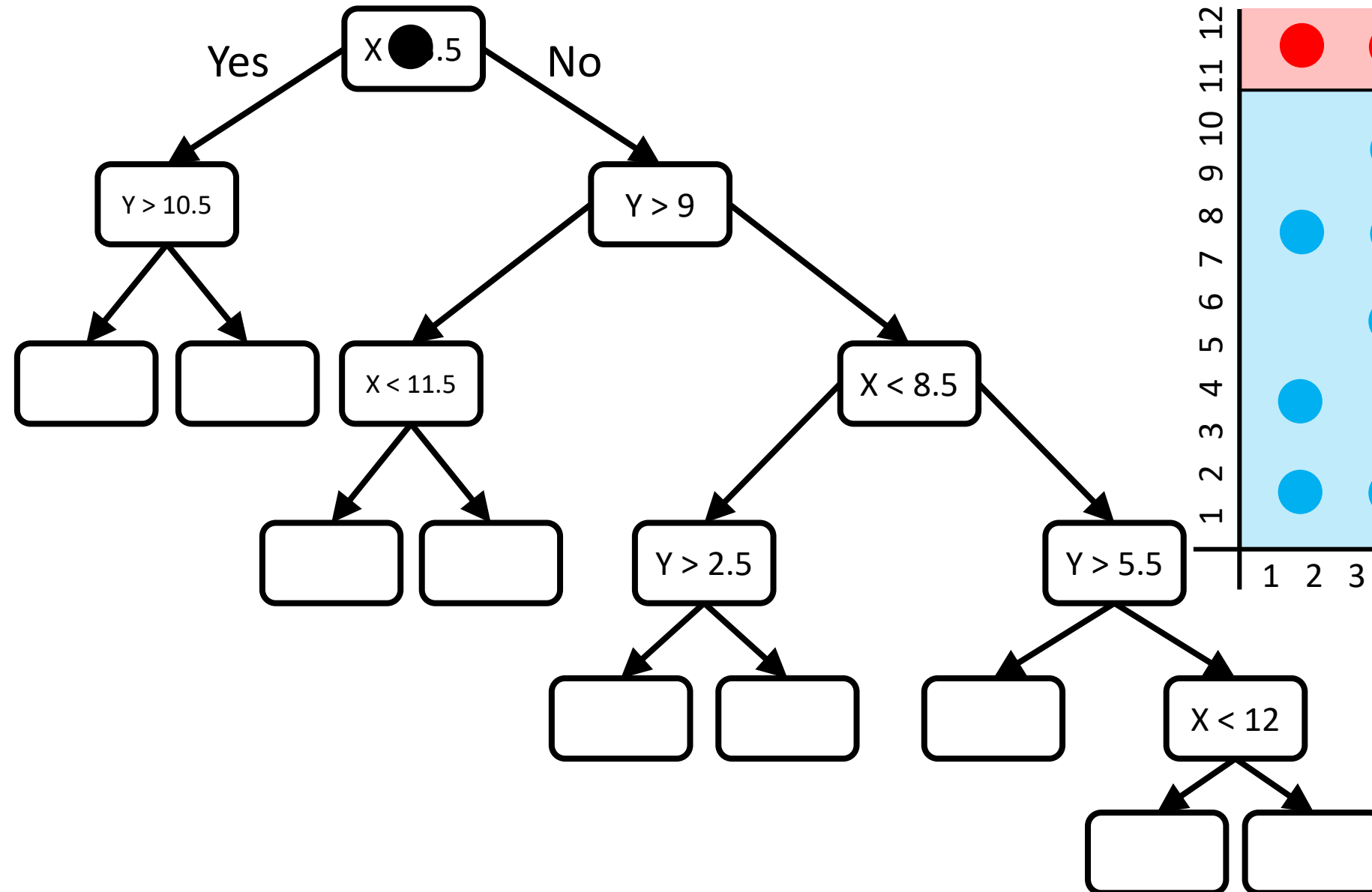
# Building Decision Trees

44



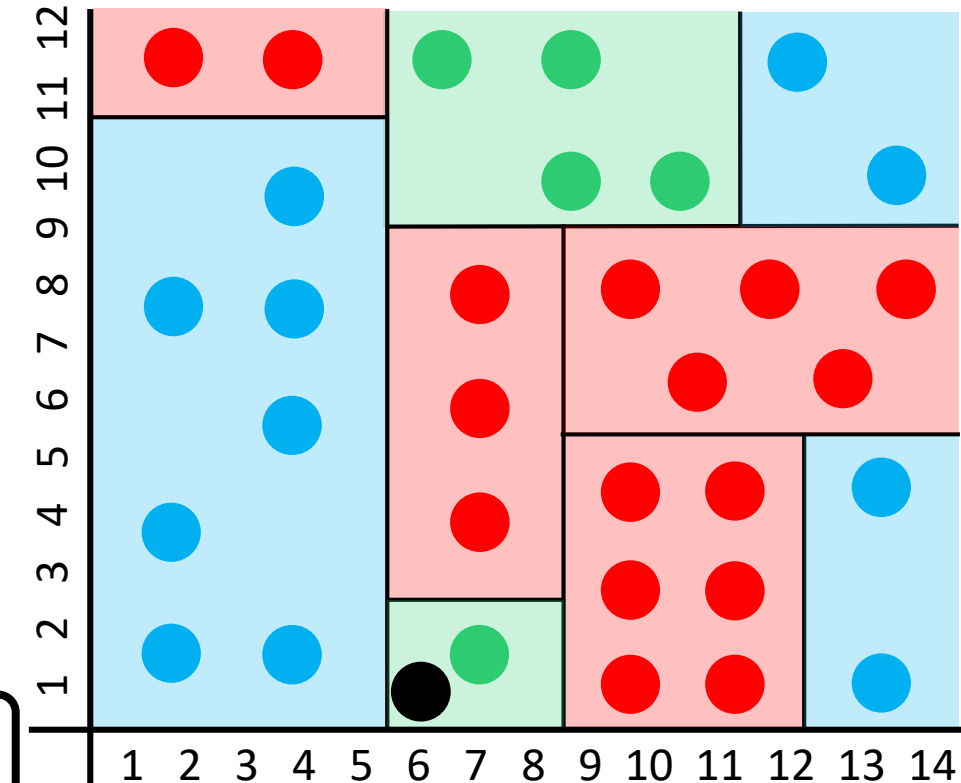
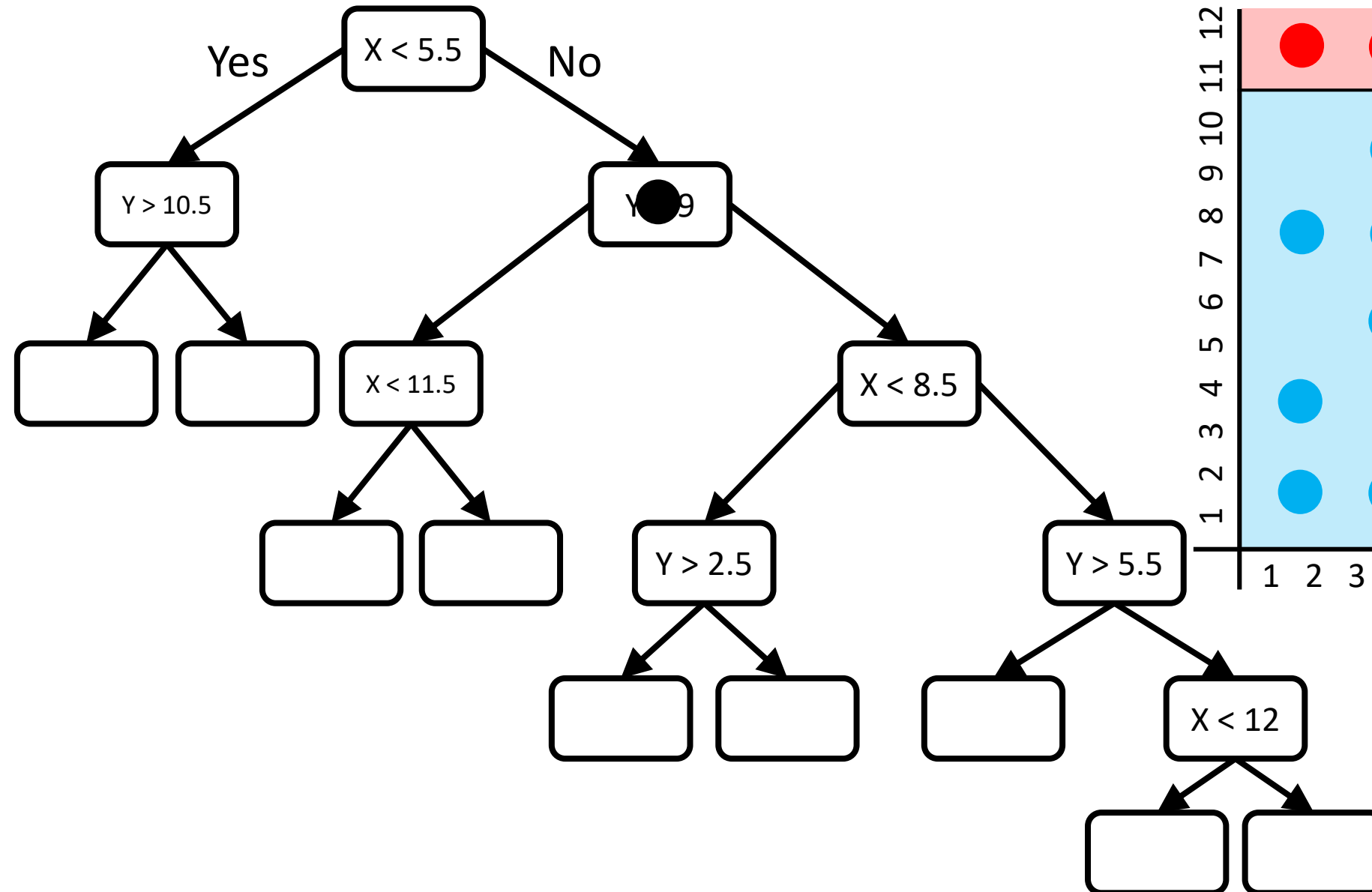
# Building Decision Trees

44



# Building Decision Trees

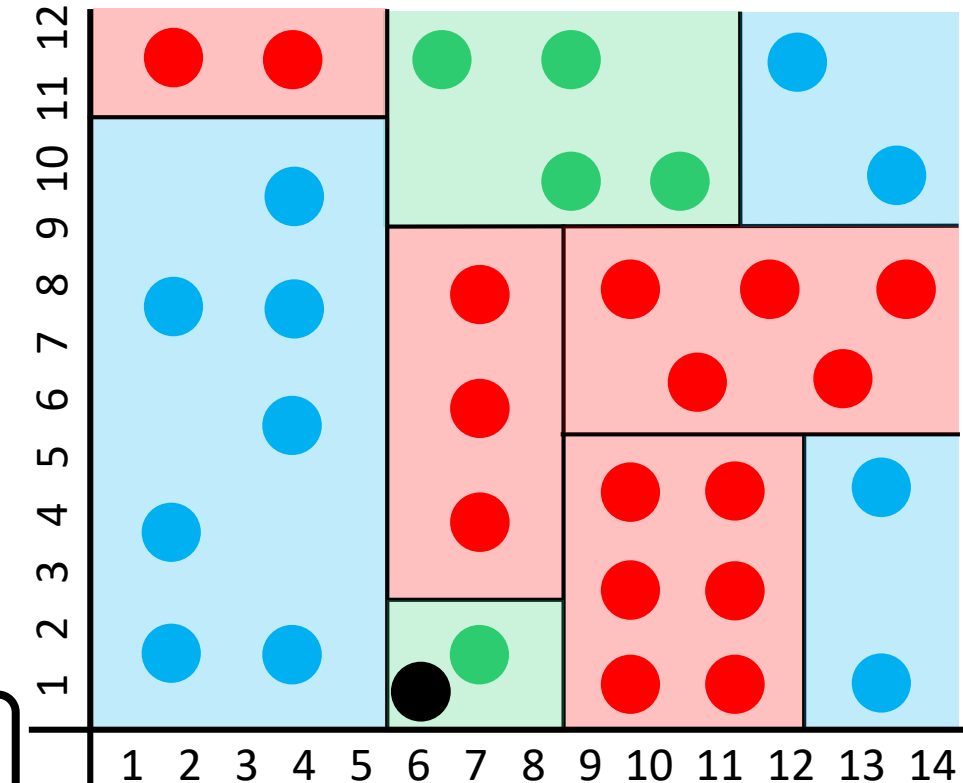
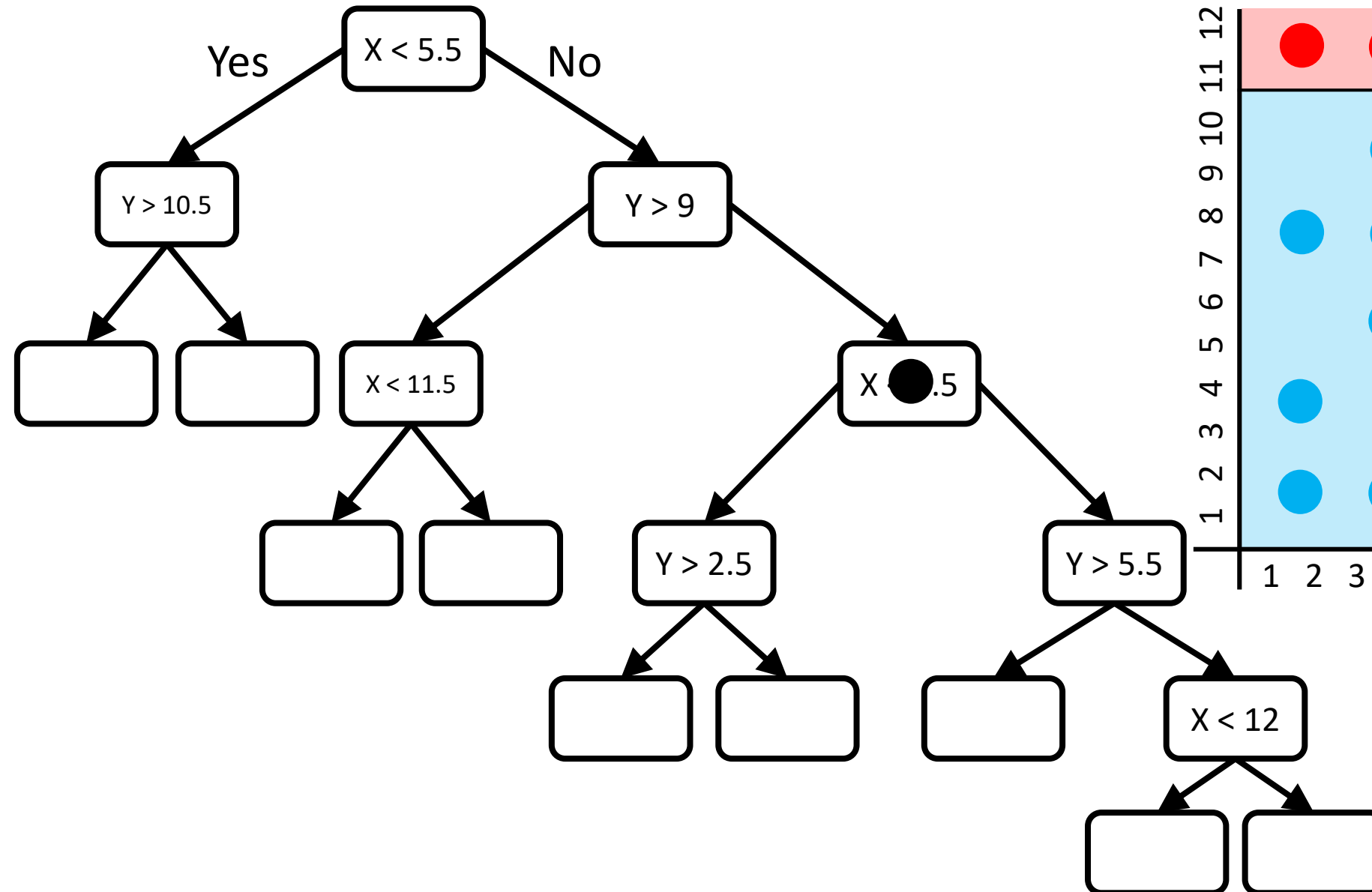
44





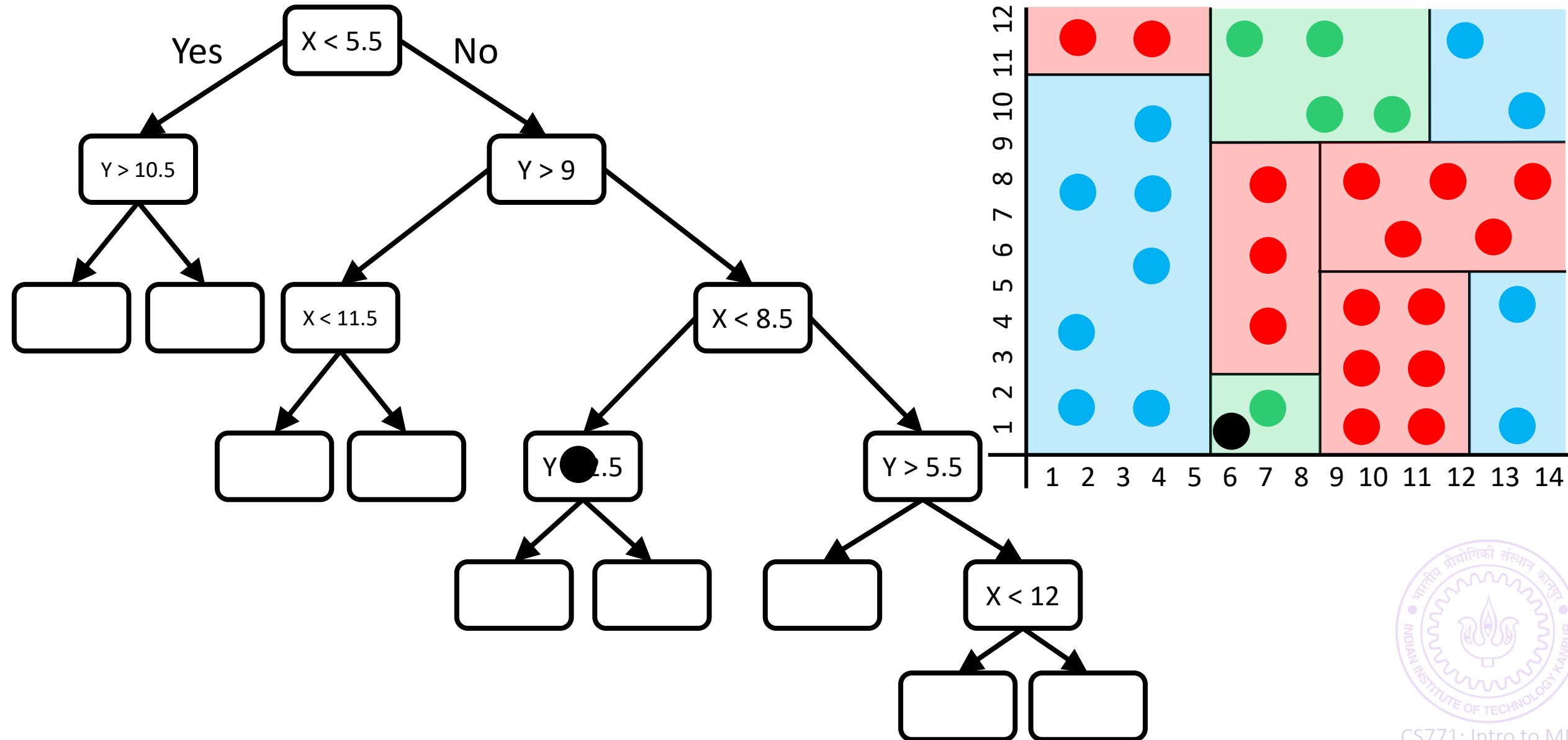
# Building Decision Trees

44



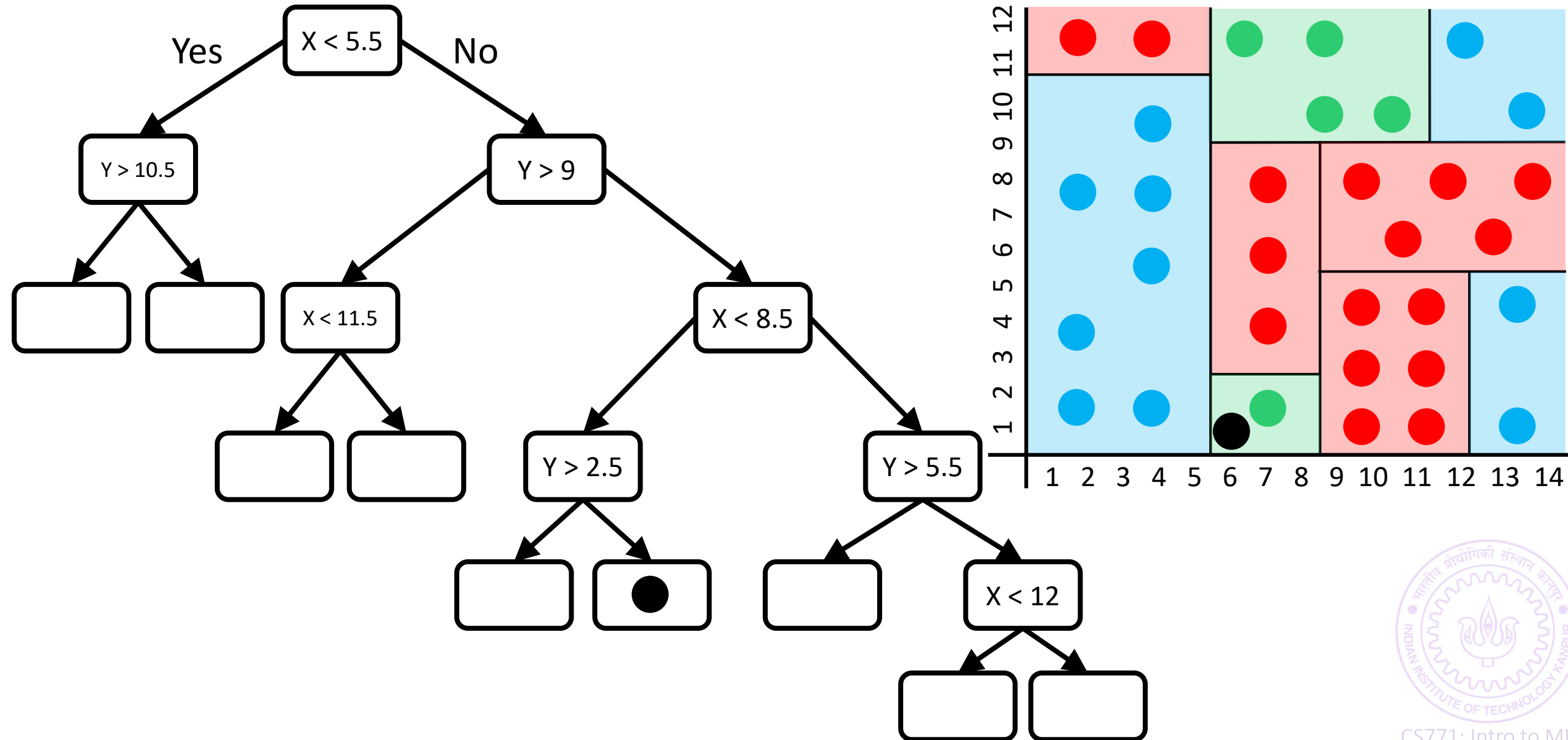
# Building Decision Trees

44



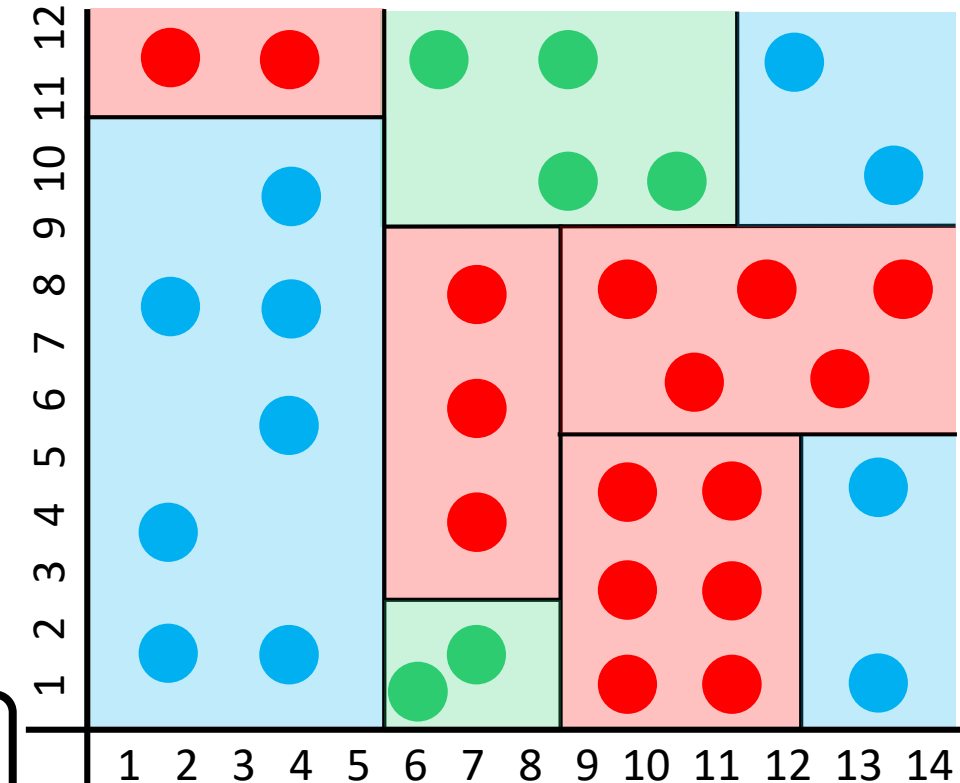
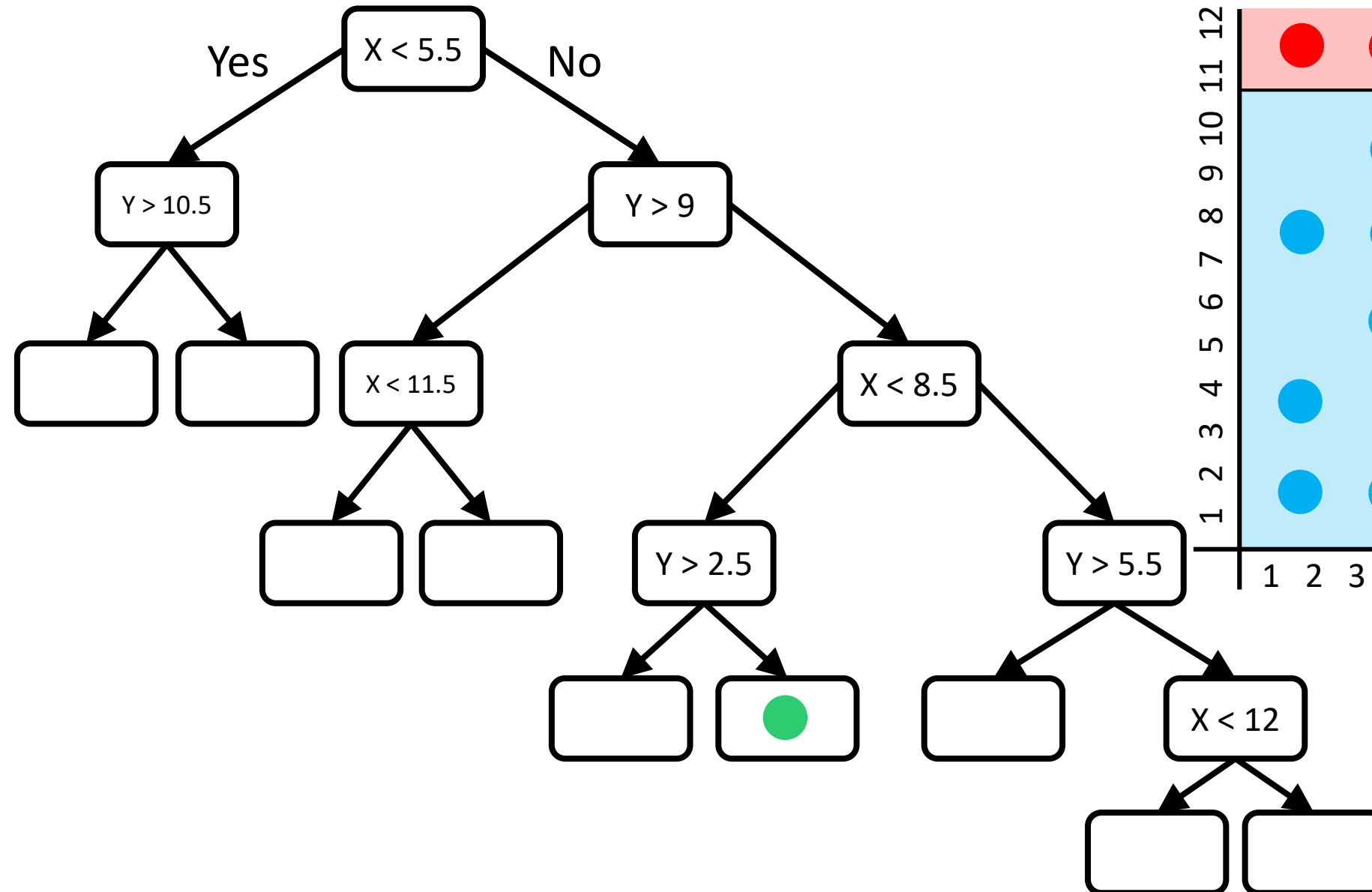
# Building Decision Trees

44



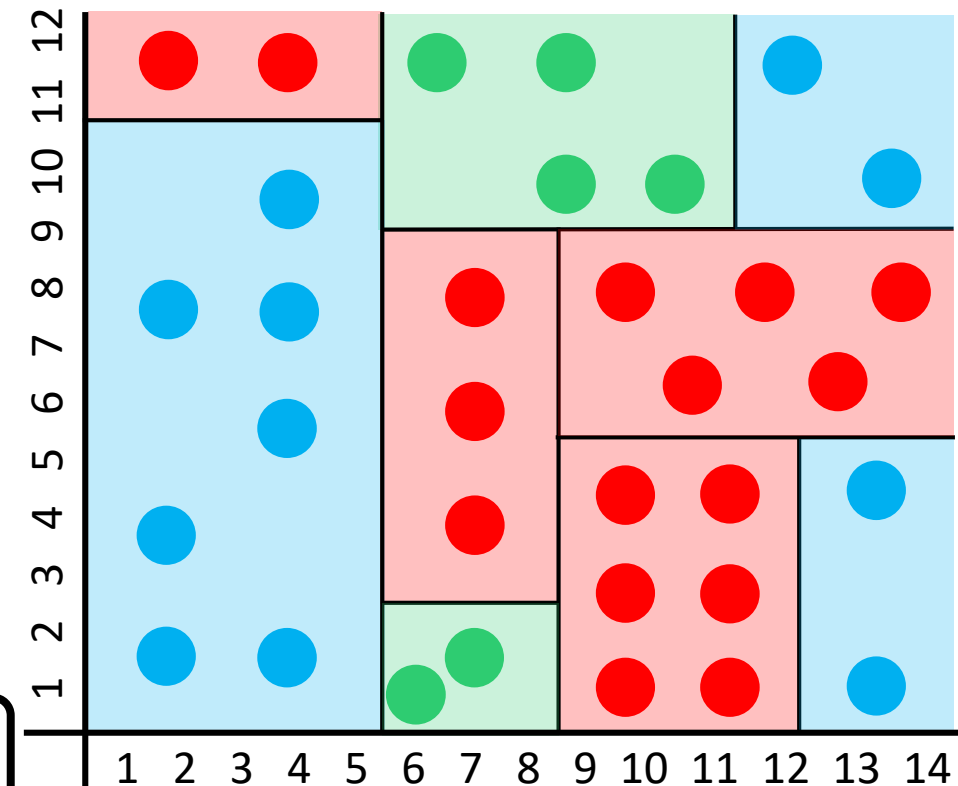
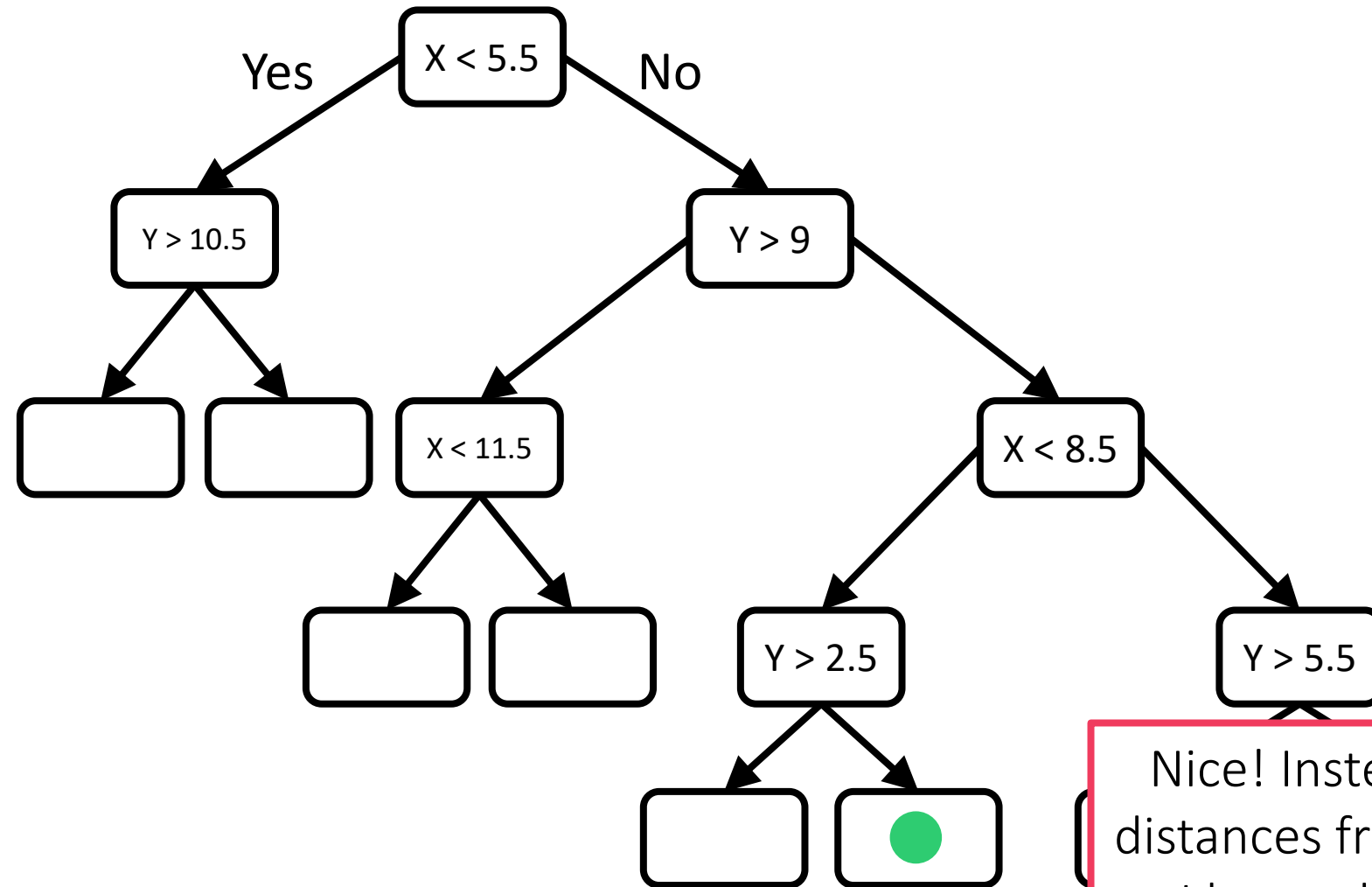
# Building Decision Trees

44

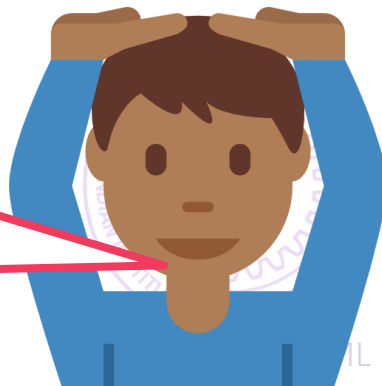


# Building Decision Trees

44

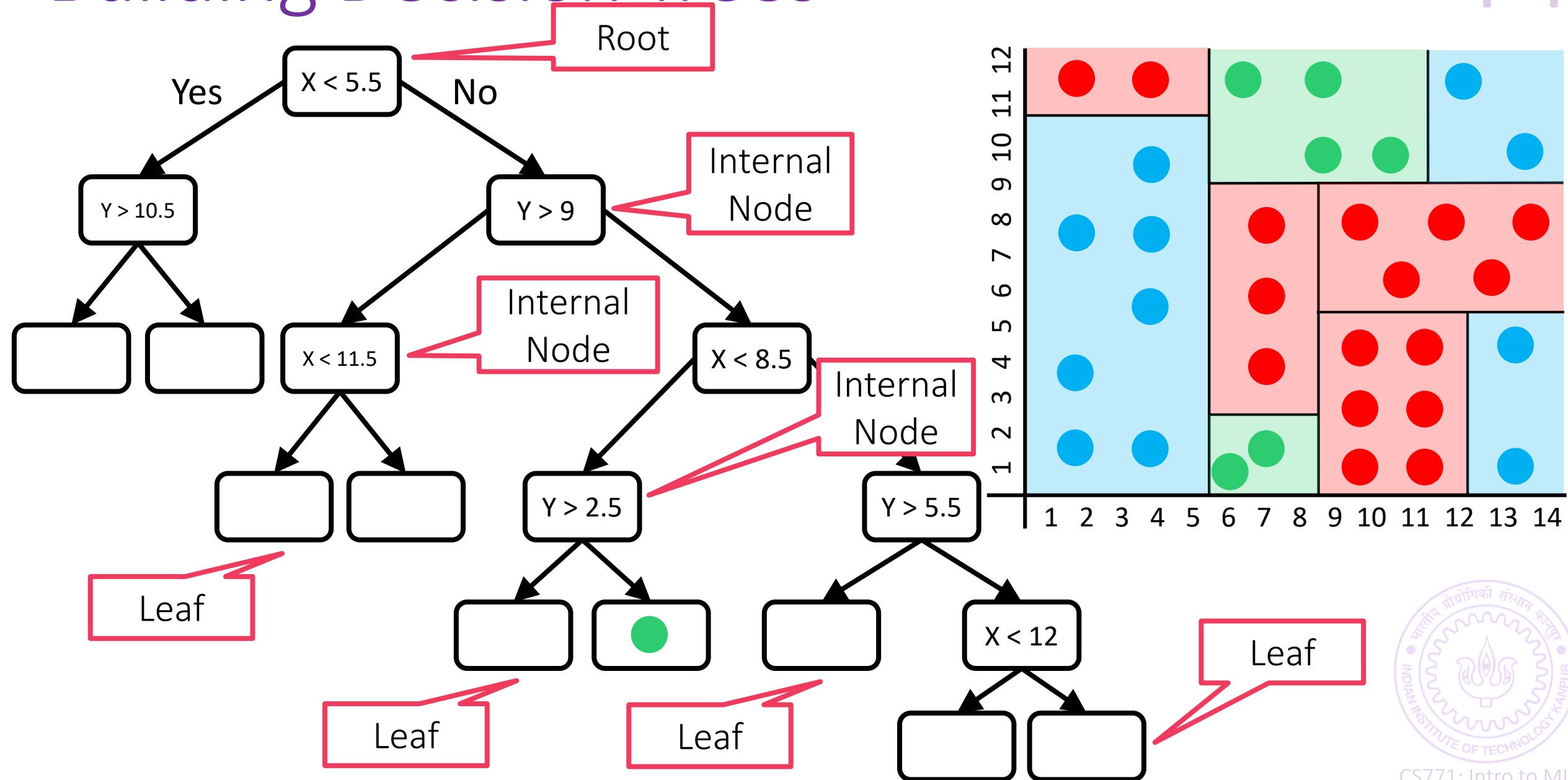


Nice! Instead of calculating distances from 32 data points, I located my partition by evaluating just 4 inequalities!!



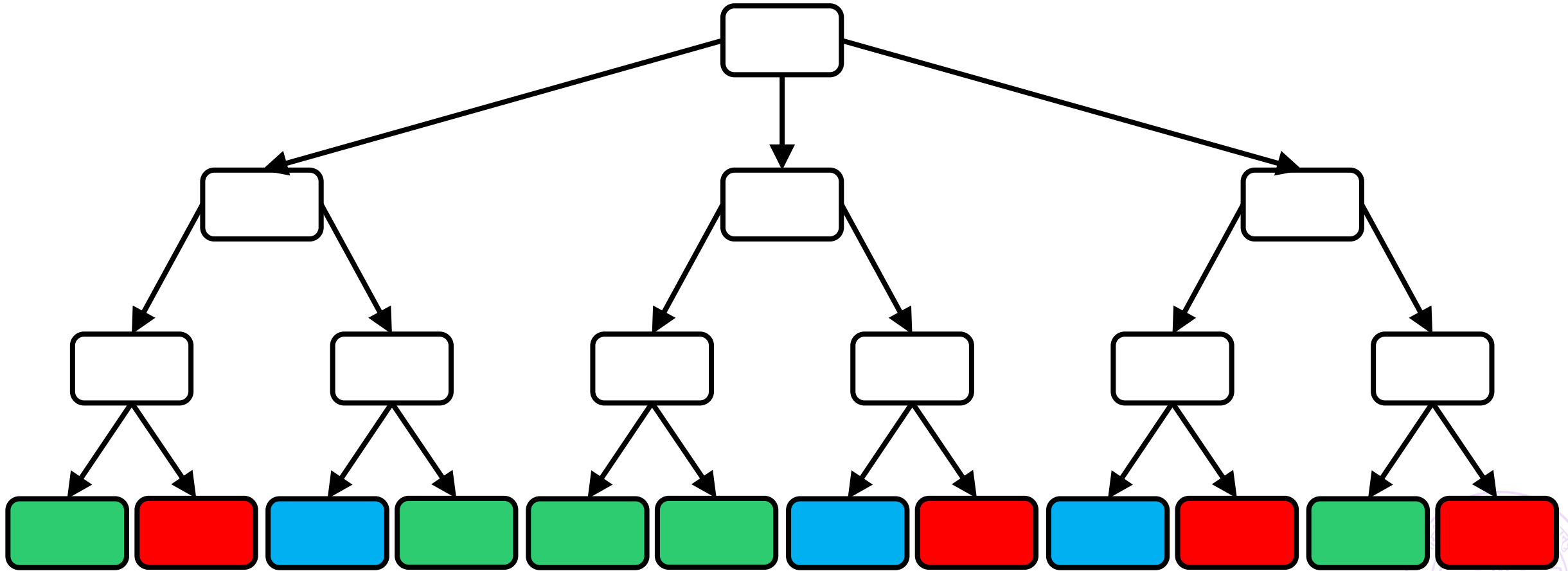
# Building Decision Trees

44



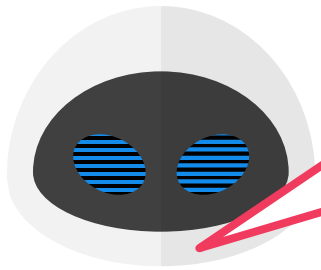
# Decision Trees – all shapes and sizes

63

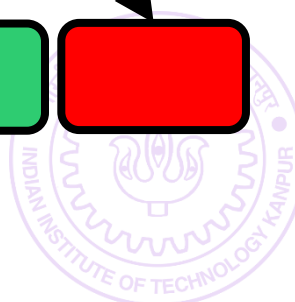
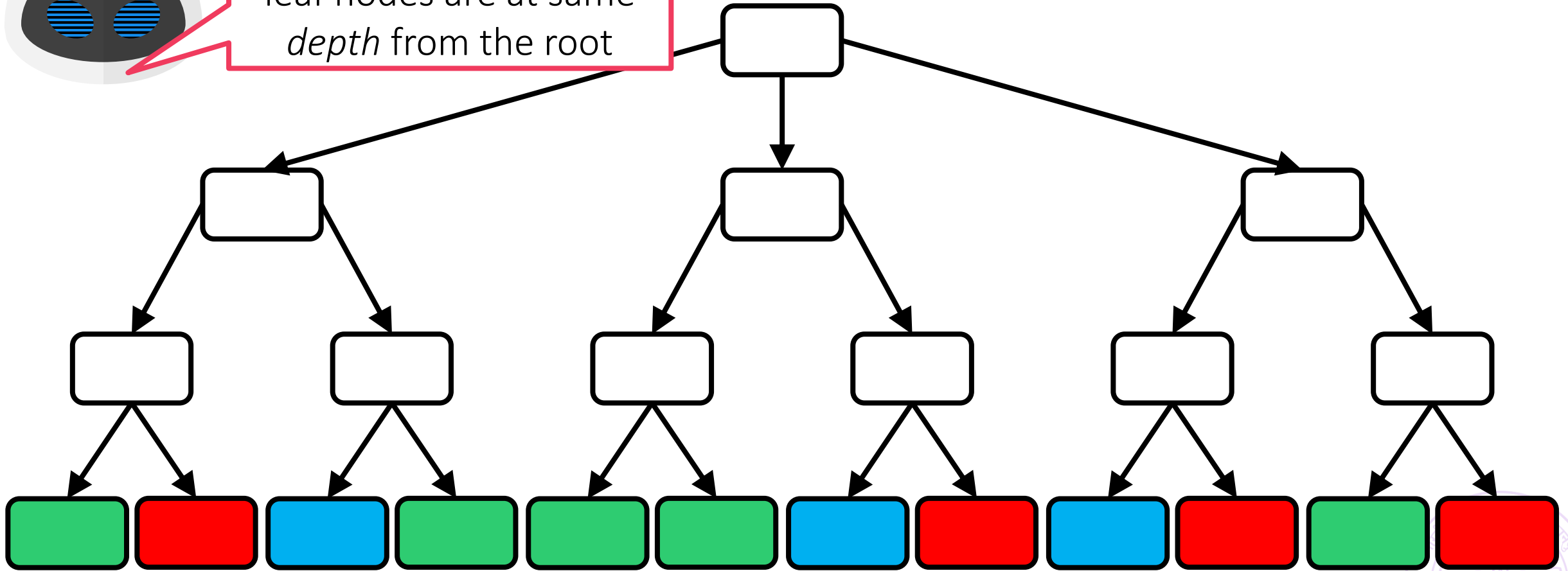


# Decision Trees – all shapes and sizes

63



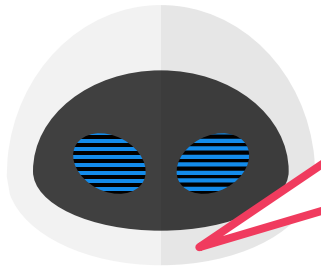
This DT is *balanced* – all leaf nodes are at same *depth* from the root





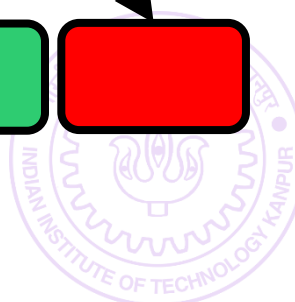
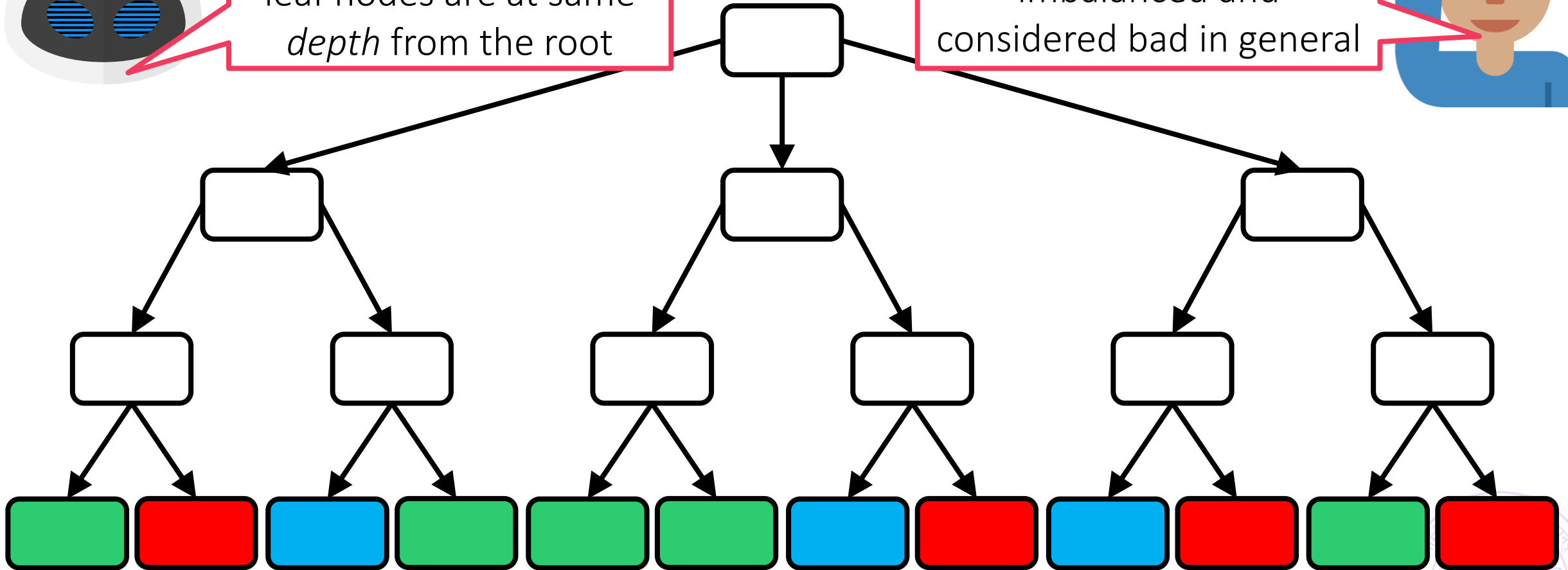
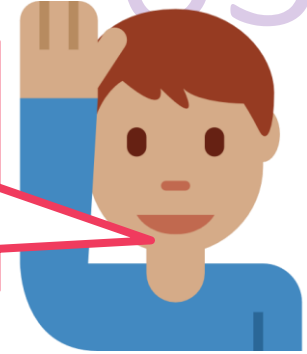
# Decision Trees – all shapes and sizes

63



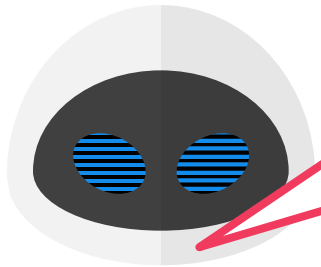
This DT is *balanced* – all leaf nodes are at same *depth* from the root

The previous DT was very imbalanced and considered bad in general

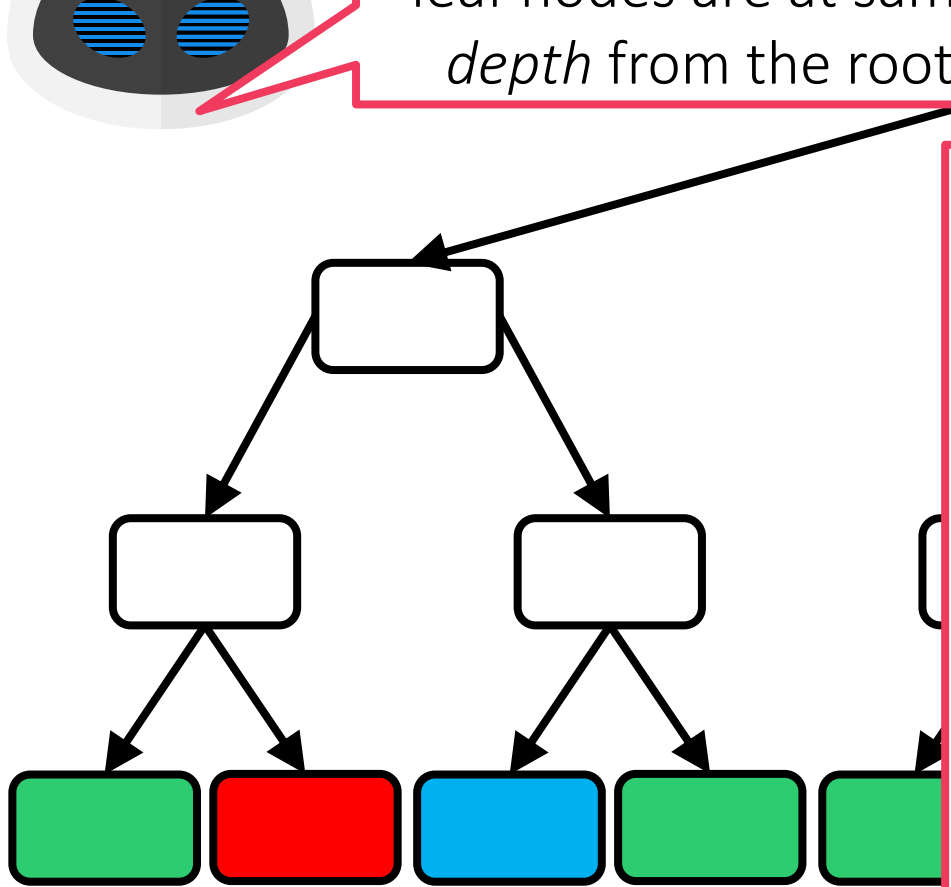


# Decision Trees – all shapes and sizes

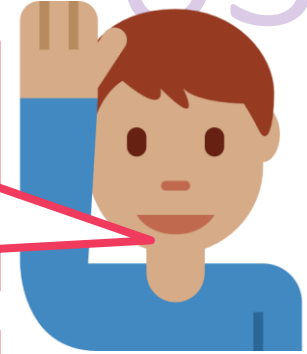
63



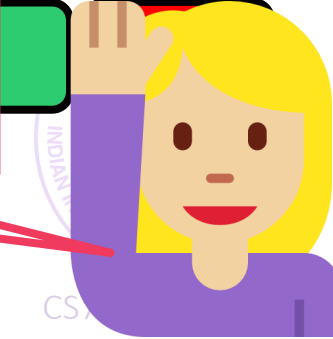
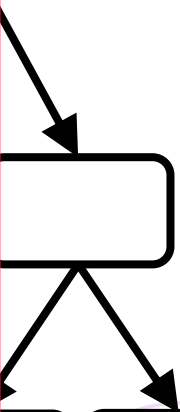
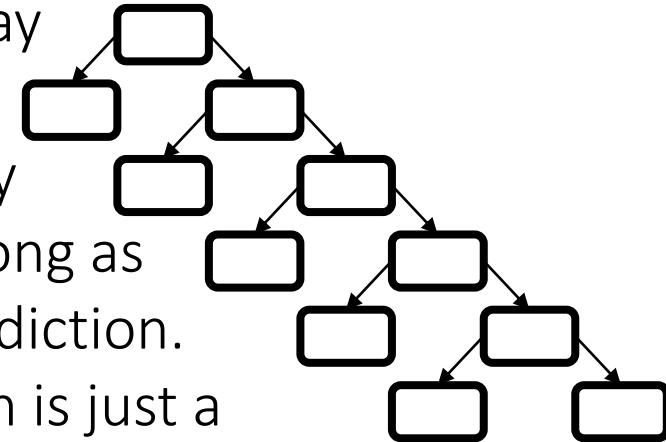
This DT is *balanced* – all leaf nodes are at same *depth* from the root



The previous DT was very imbalanced and considered bad in general

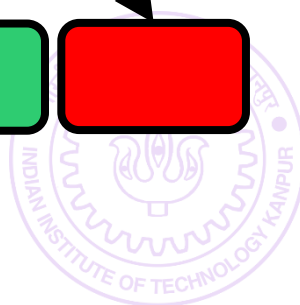
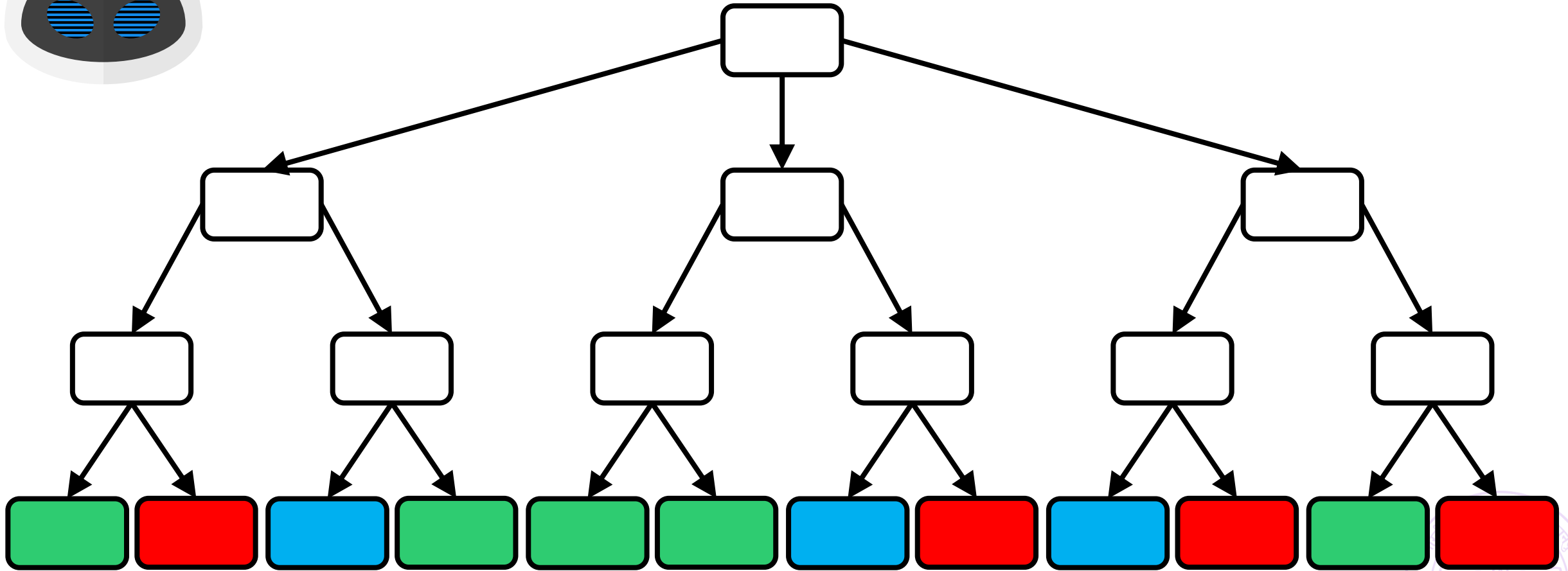
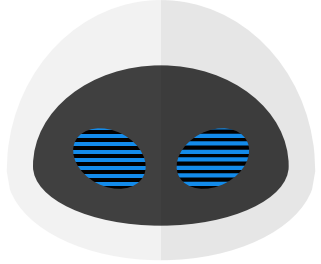


Imbalanced DTs may offer very poor Prediction accuracy as well as take as long as kNN to make a prediction. Imagine a DT which is just a chain of nodes. With  $n$  data points, there would be  $\mathcal{O}(n)$  chain: some predictions will take  $\mathcal{O}(n)$  time ☹️. With a balanced DT, every prediction takes at most  $\mathcal{O}(\log n)$  time 😊😊



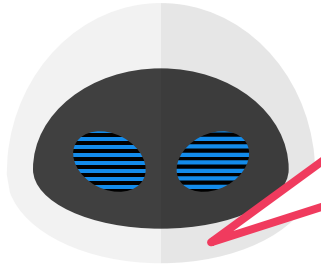
# Decision Trees – all shapes and sizes

63

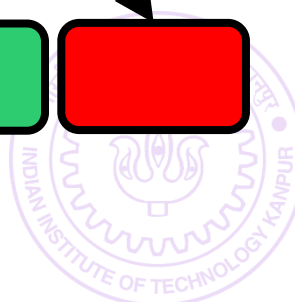
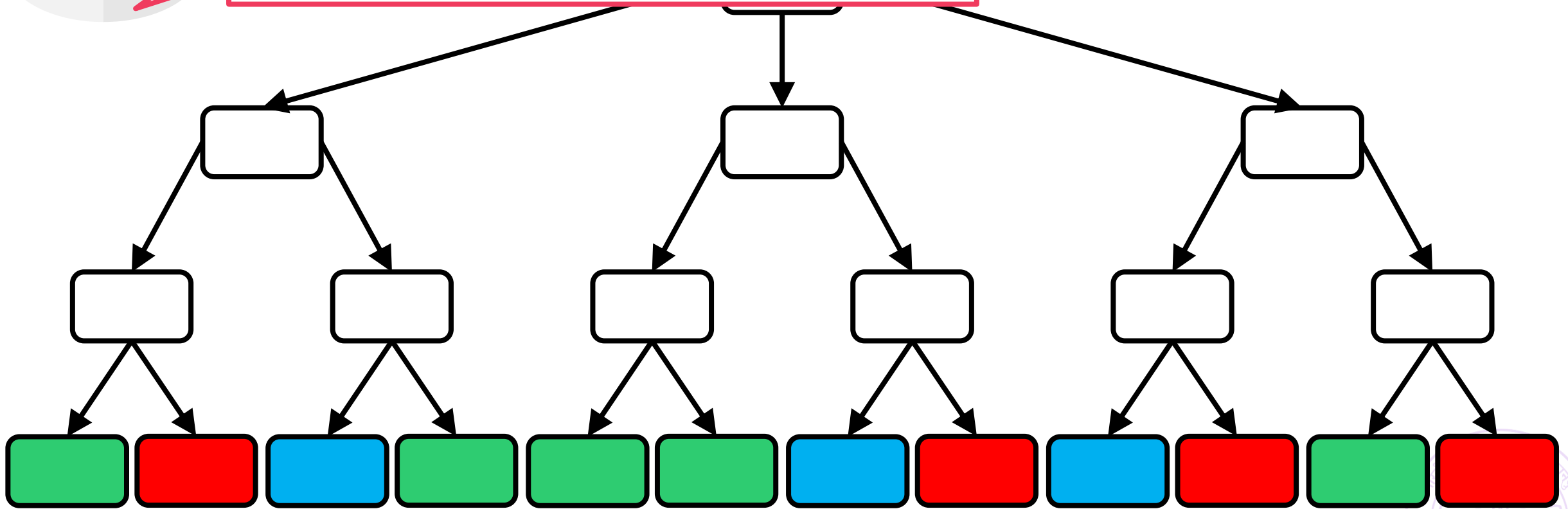


# Decision Trees – all shapes and sizes

63

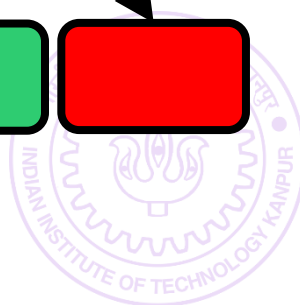
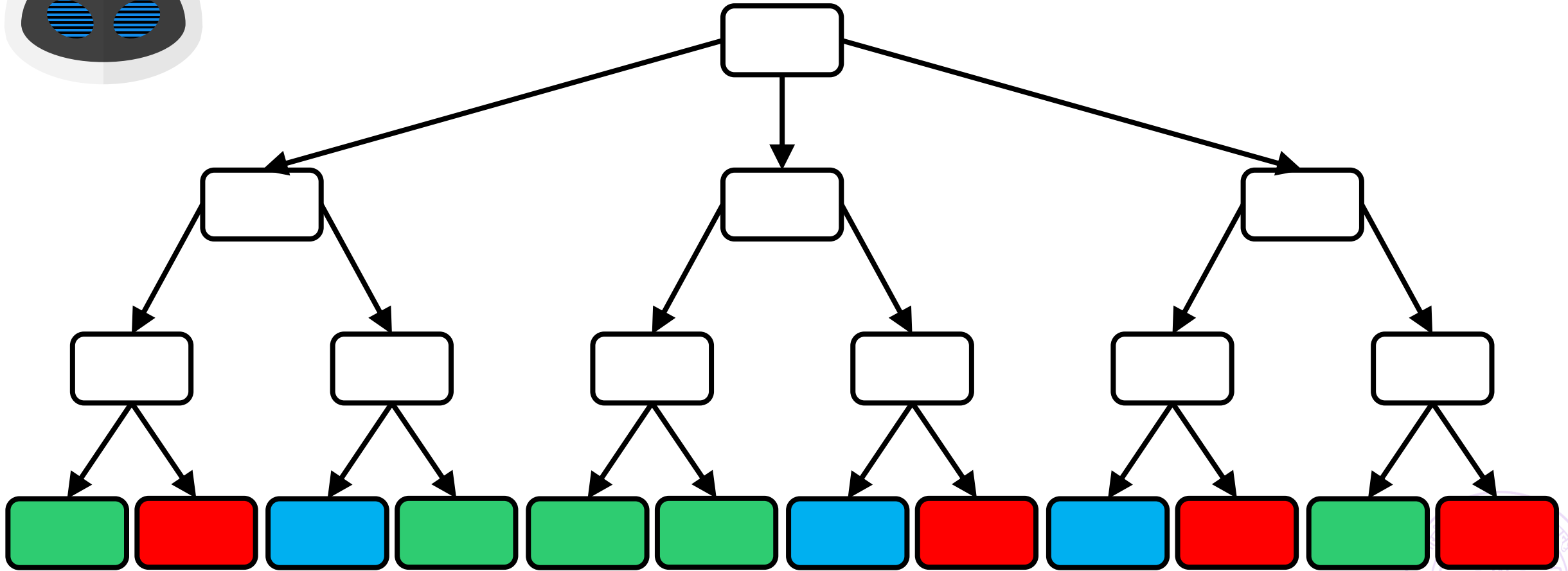
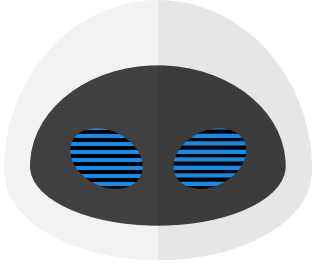


May prune the tree to make it more shallow as well. Also possible to have DT with more than 2 children per internal node



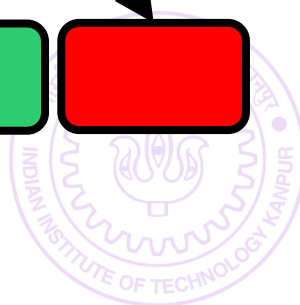
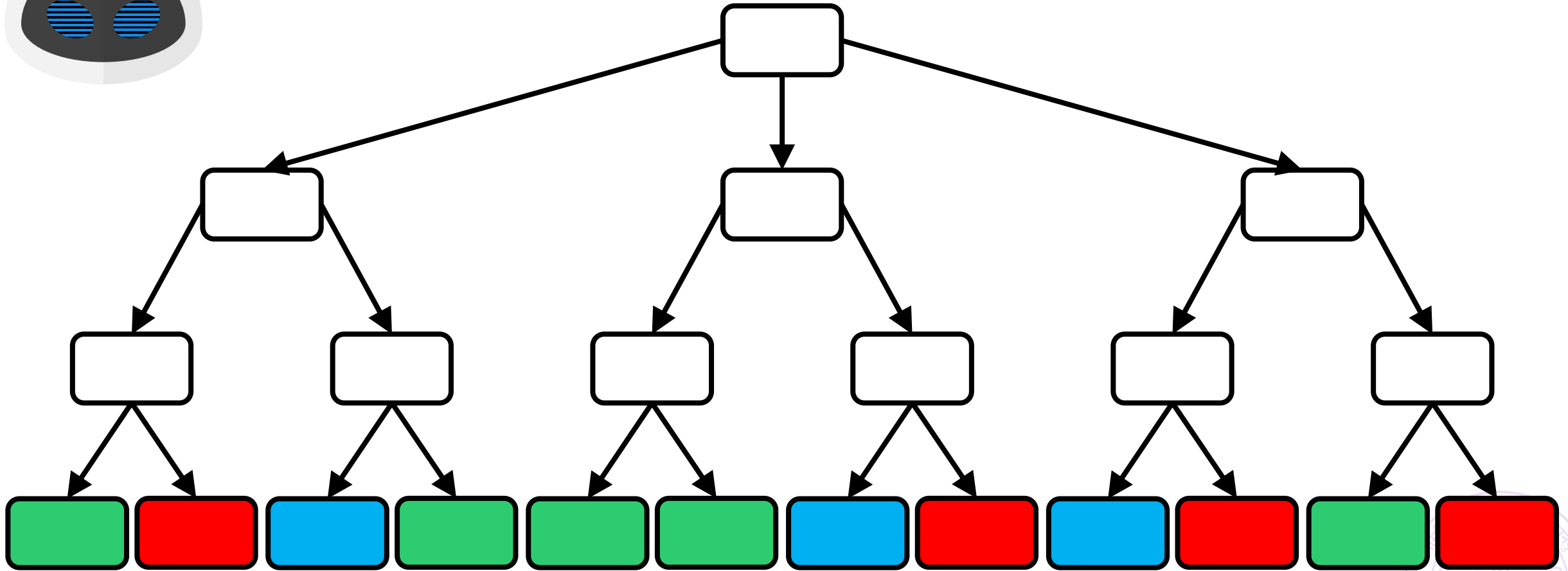
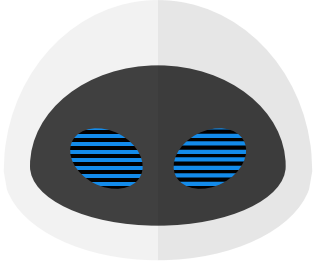
# Decision Trees – all shapes and sizes

63



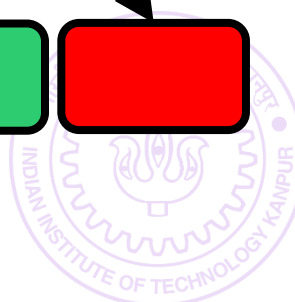
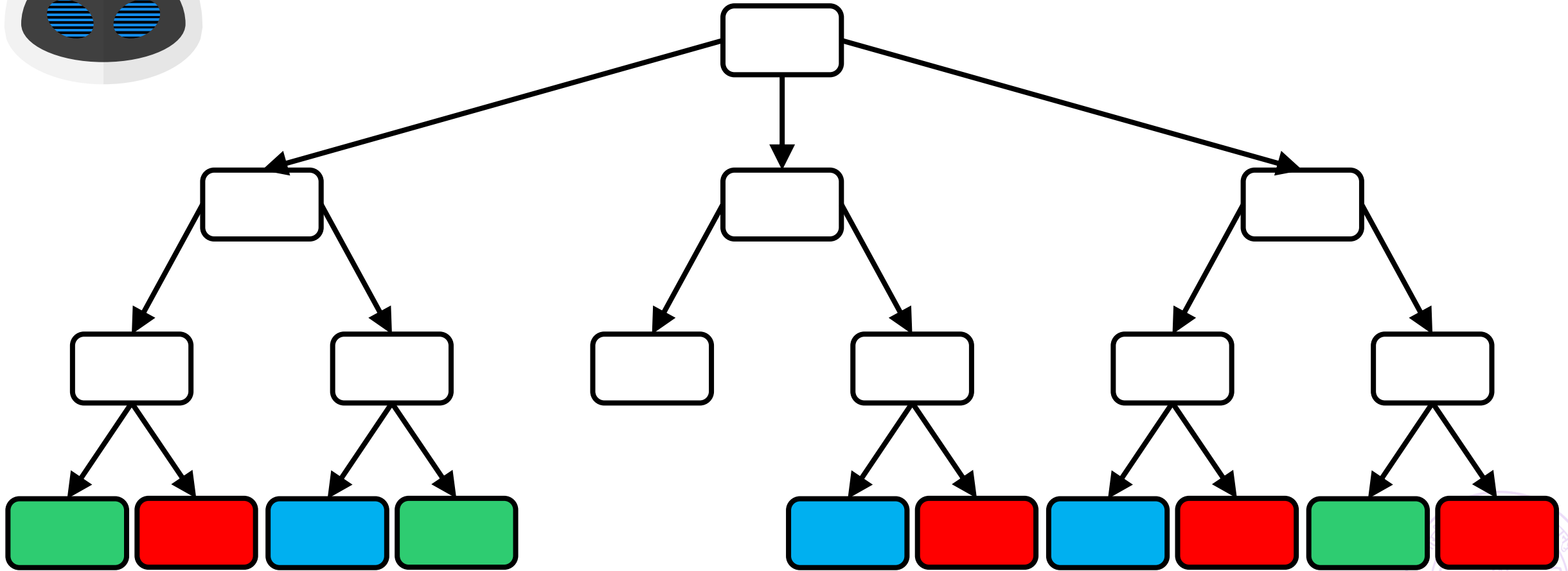
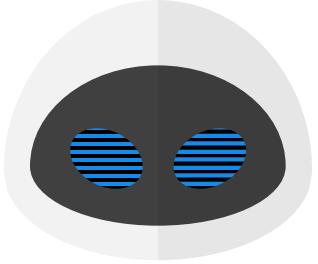
# Decision Trees – all shapes and sizes

63



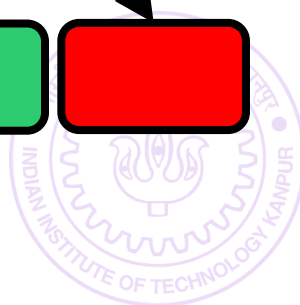
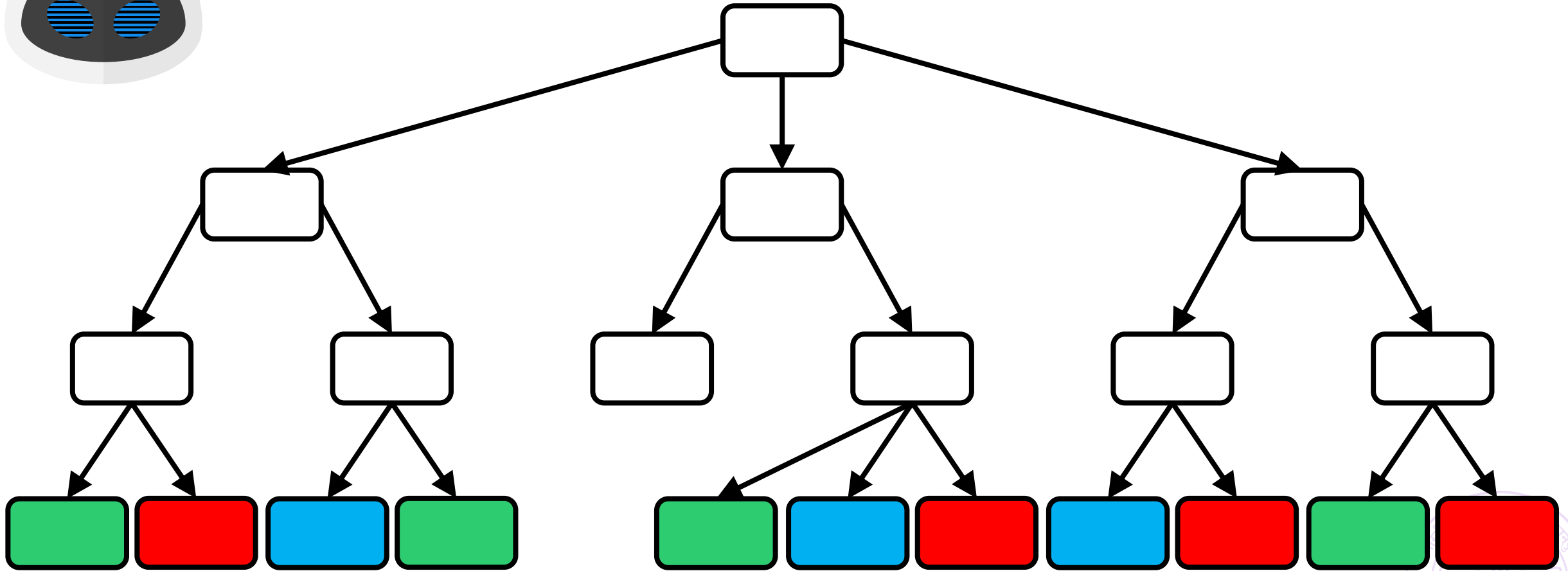
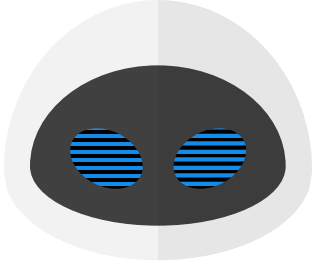
# Decision Trees – all shapes and sizes

63



# Decision Trees – all shapes and sizes

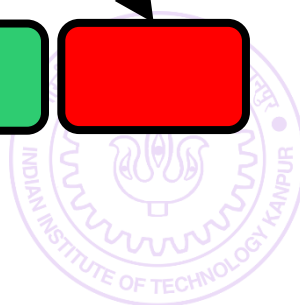
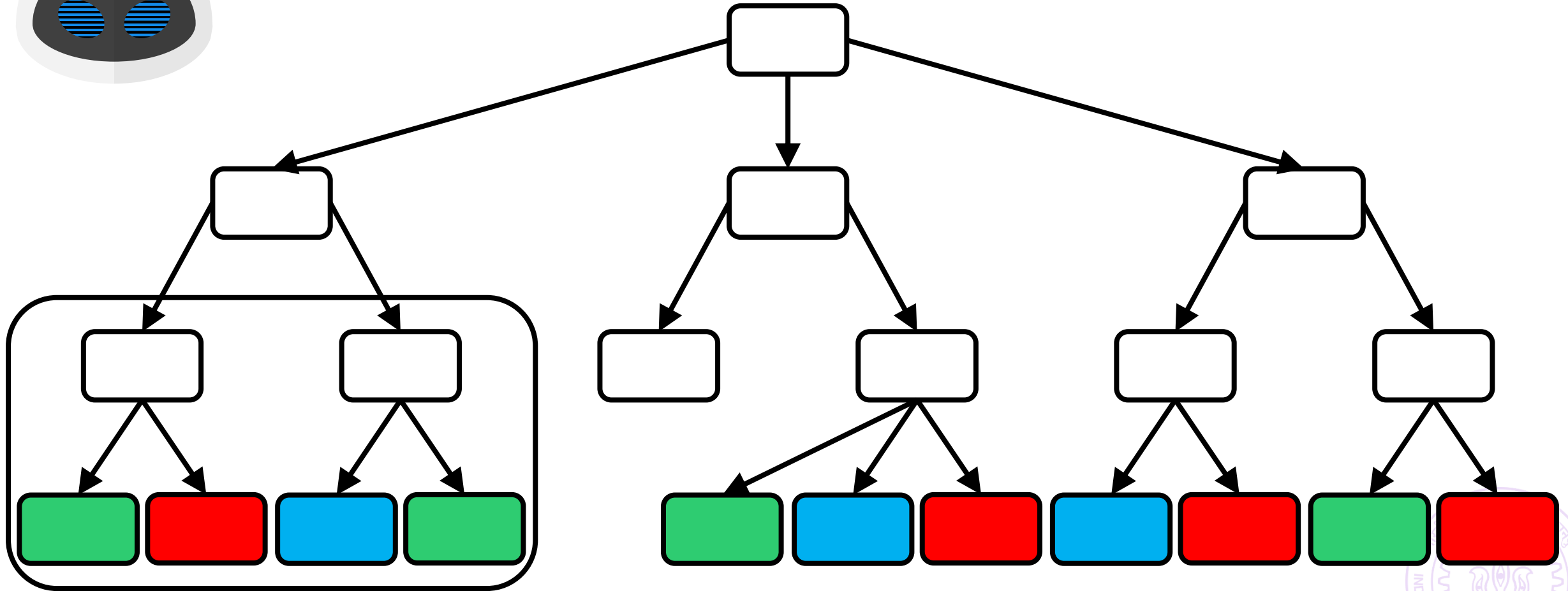
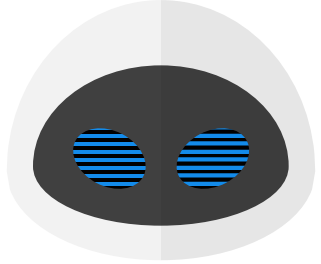
63





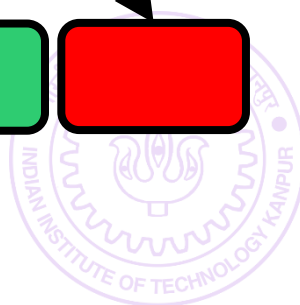
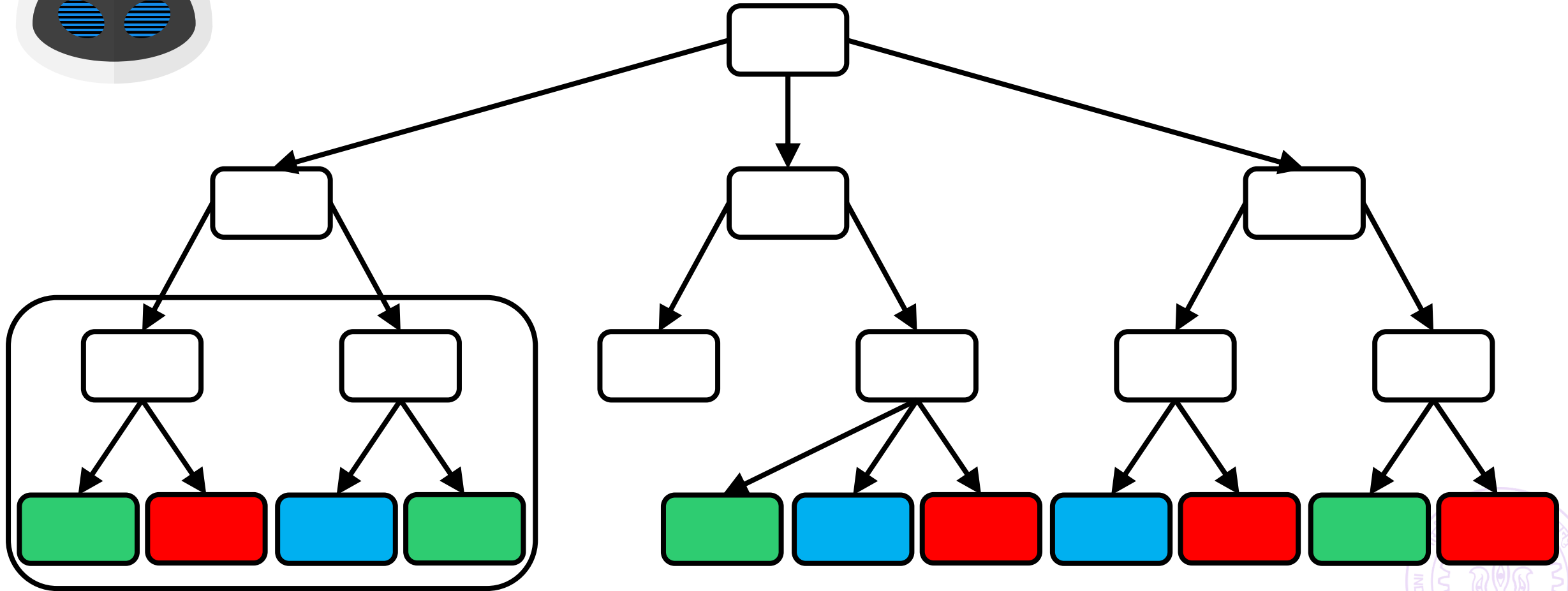
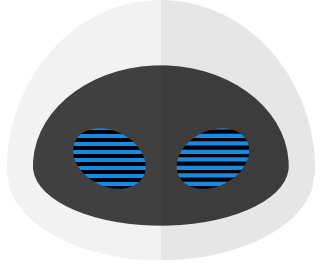
# Decision Trees – all shapes and sizes

63



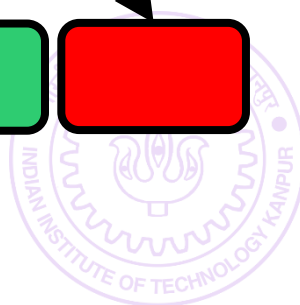
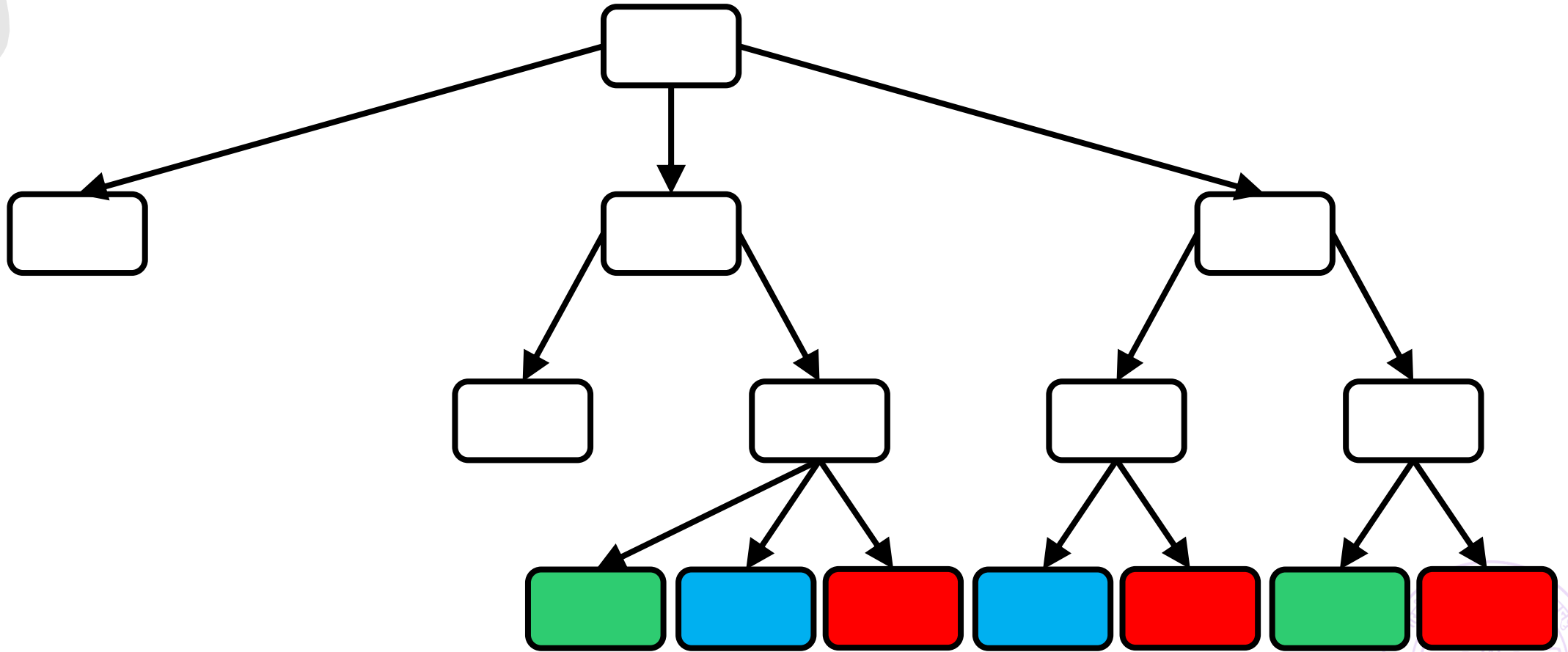
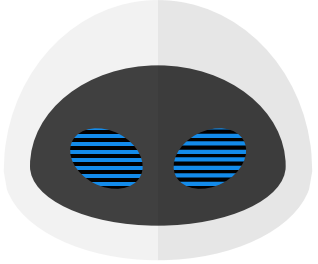
# Decision Trees – all shapes and sizes

63



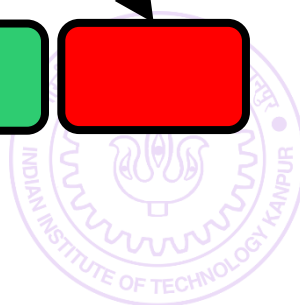
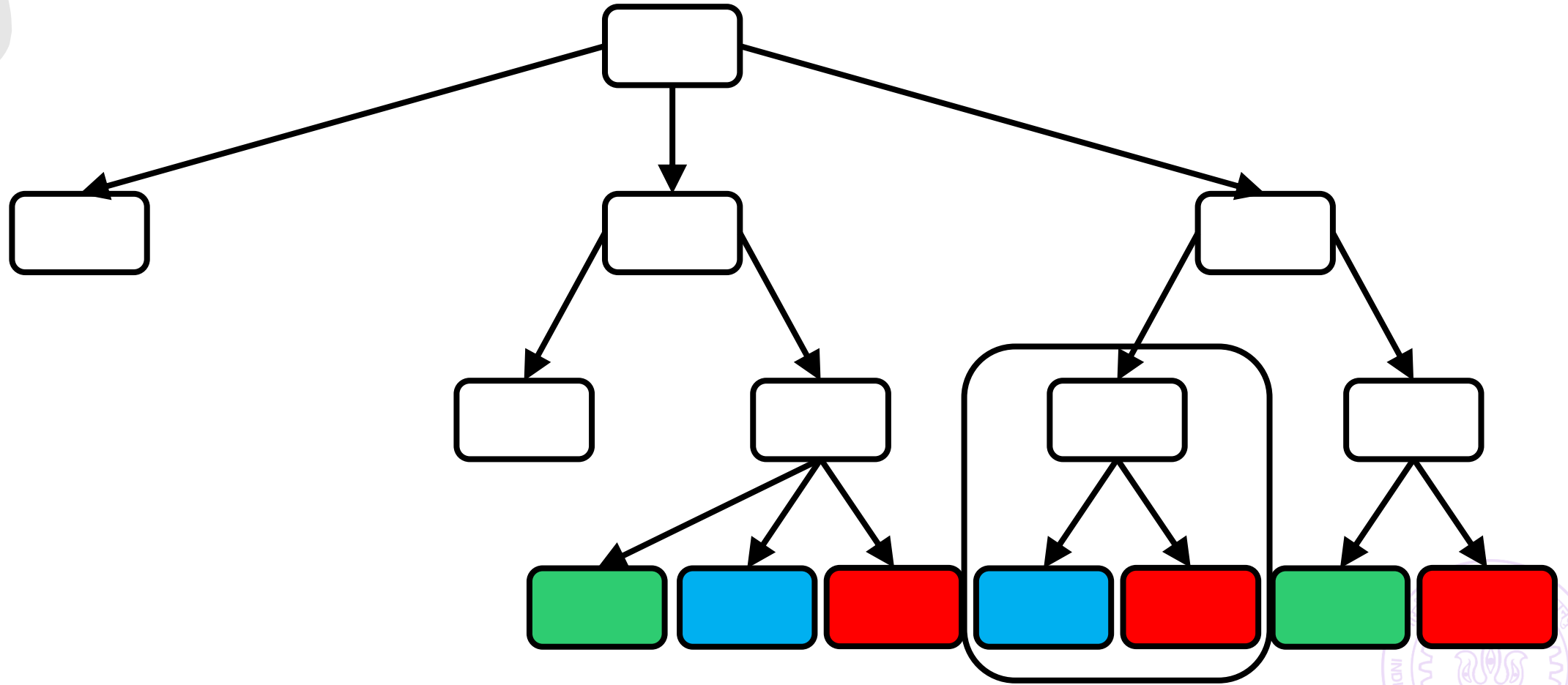
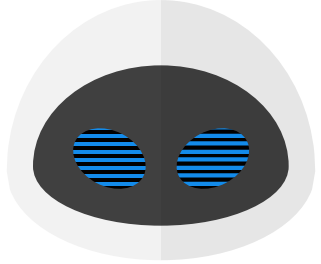
# Decision Trees – all shapes and sizes

63



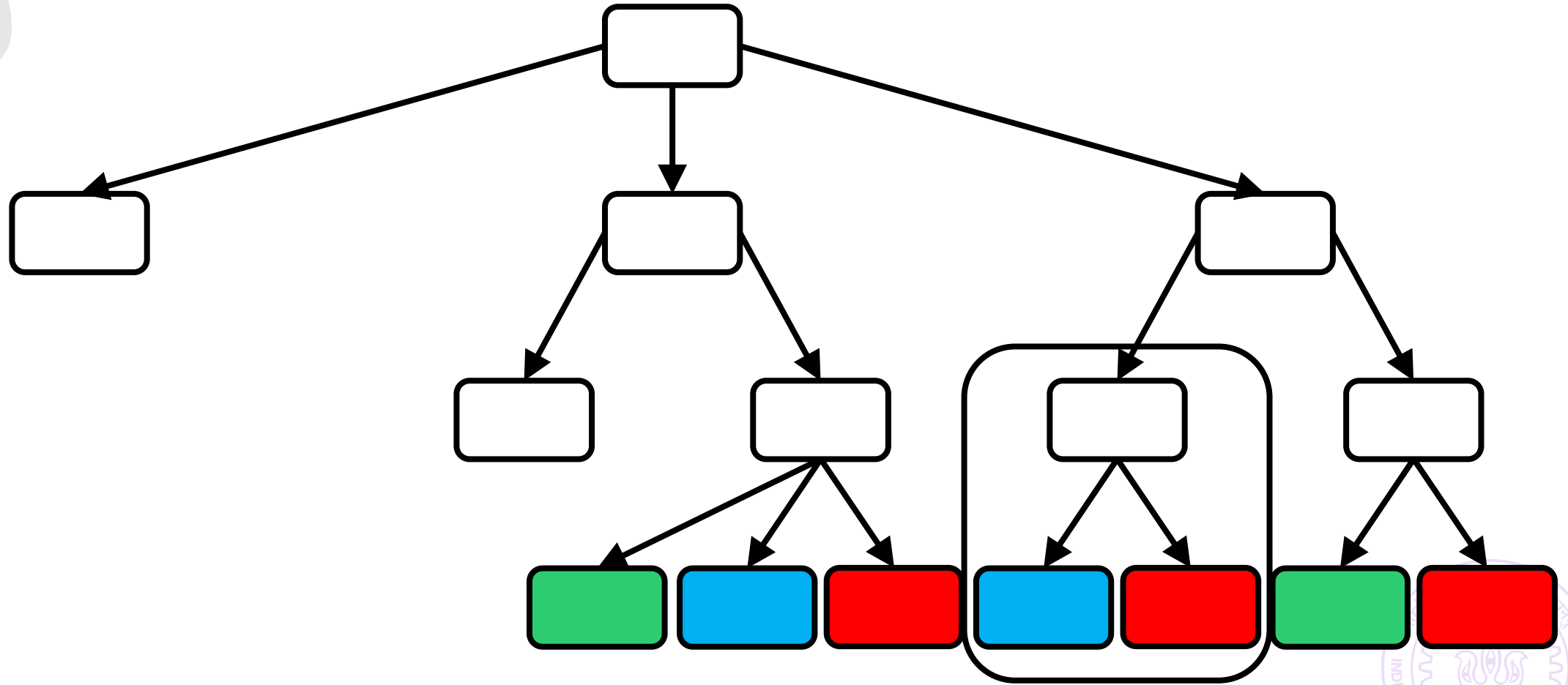
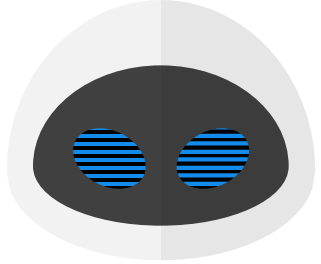
# Decision Trees – all shapes and sizes

63



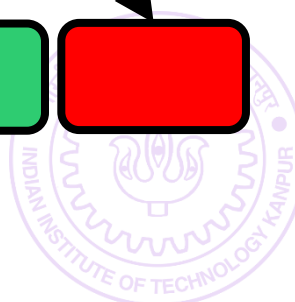
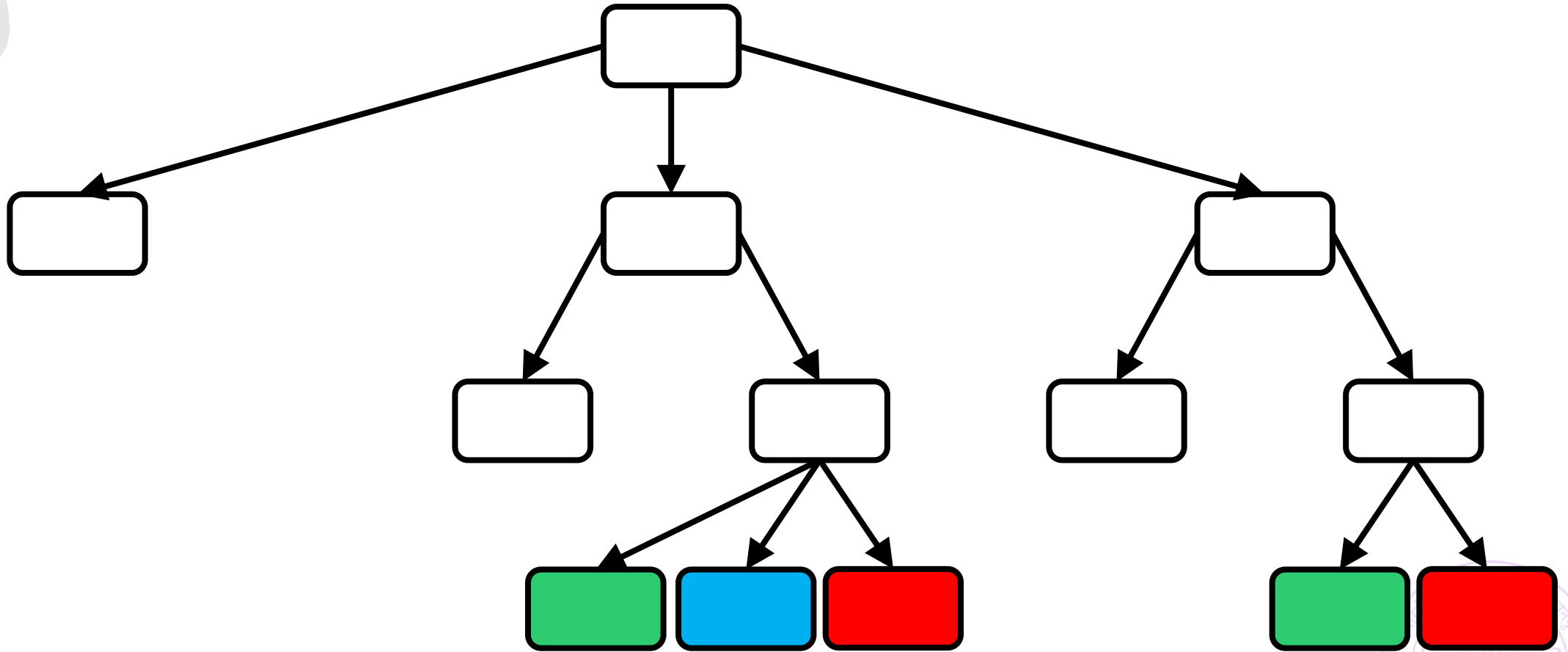
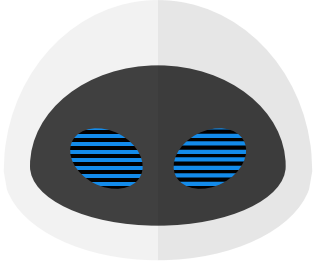
# Decision Trees – all shapes and sizes

63



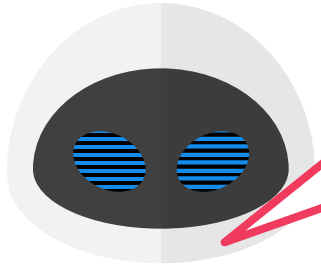
# Decision Trees – all shapes and sizes

63

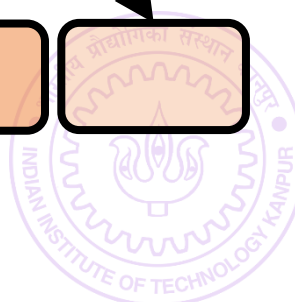
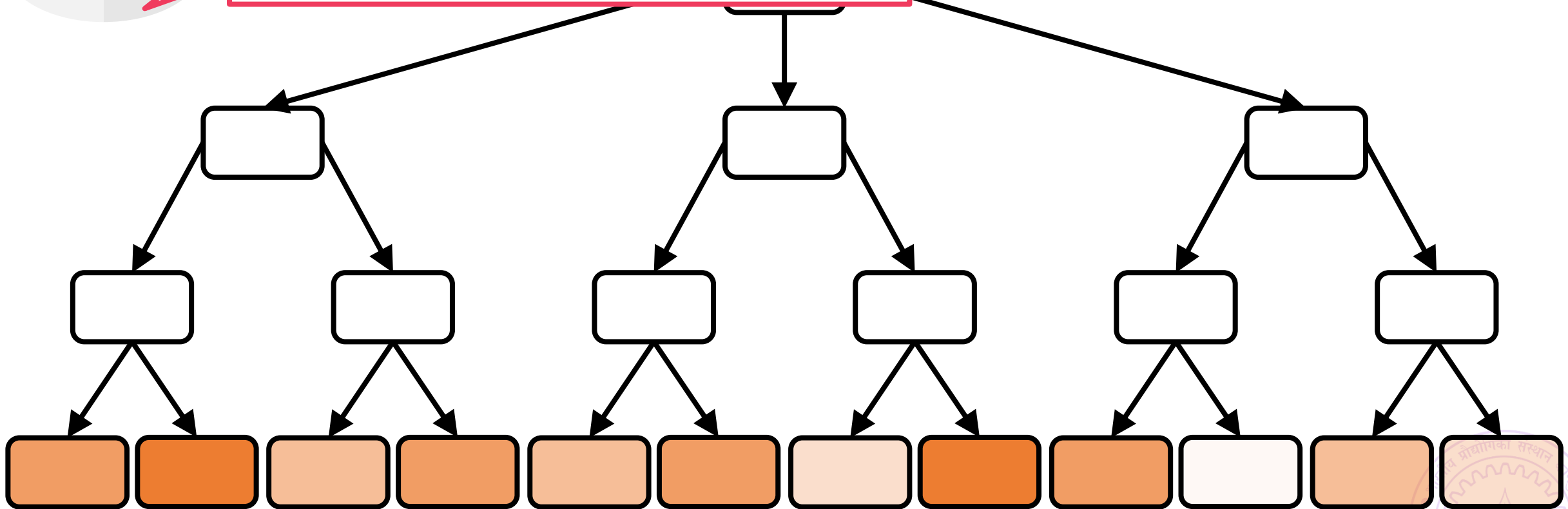


# Regression with Decision Trees

79



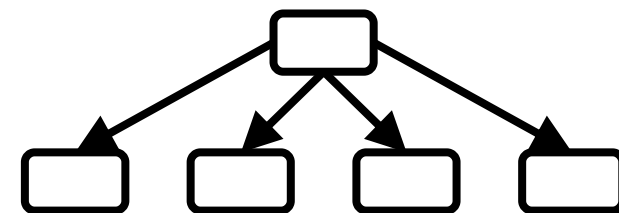
To perform real valued regression, may simply use average score at a leaf node to predict scores for test data points



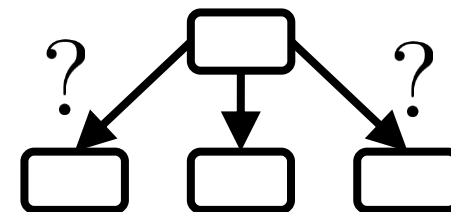
# How to learn a DT?

80

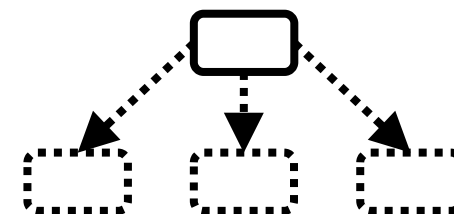
How many children should a node have?



How to send data points to children?



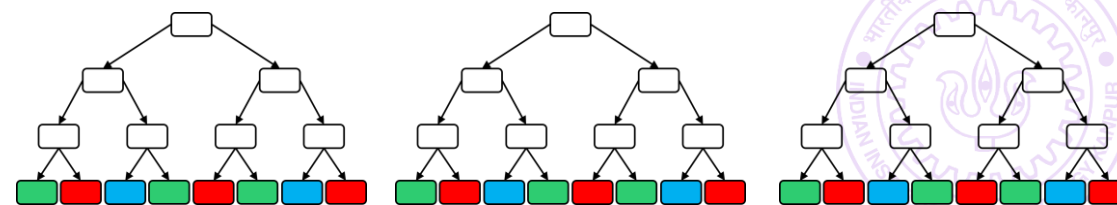
When to stop splitting and make the node a leaf?



What to do at a leaf?



How many trees to train?





# What to do at a leaf?

81

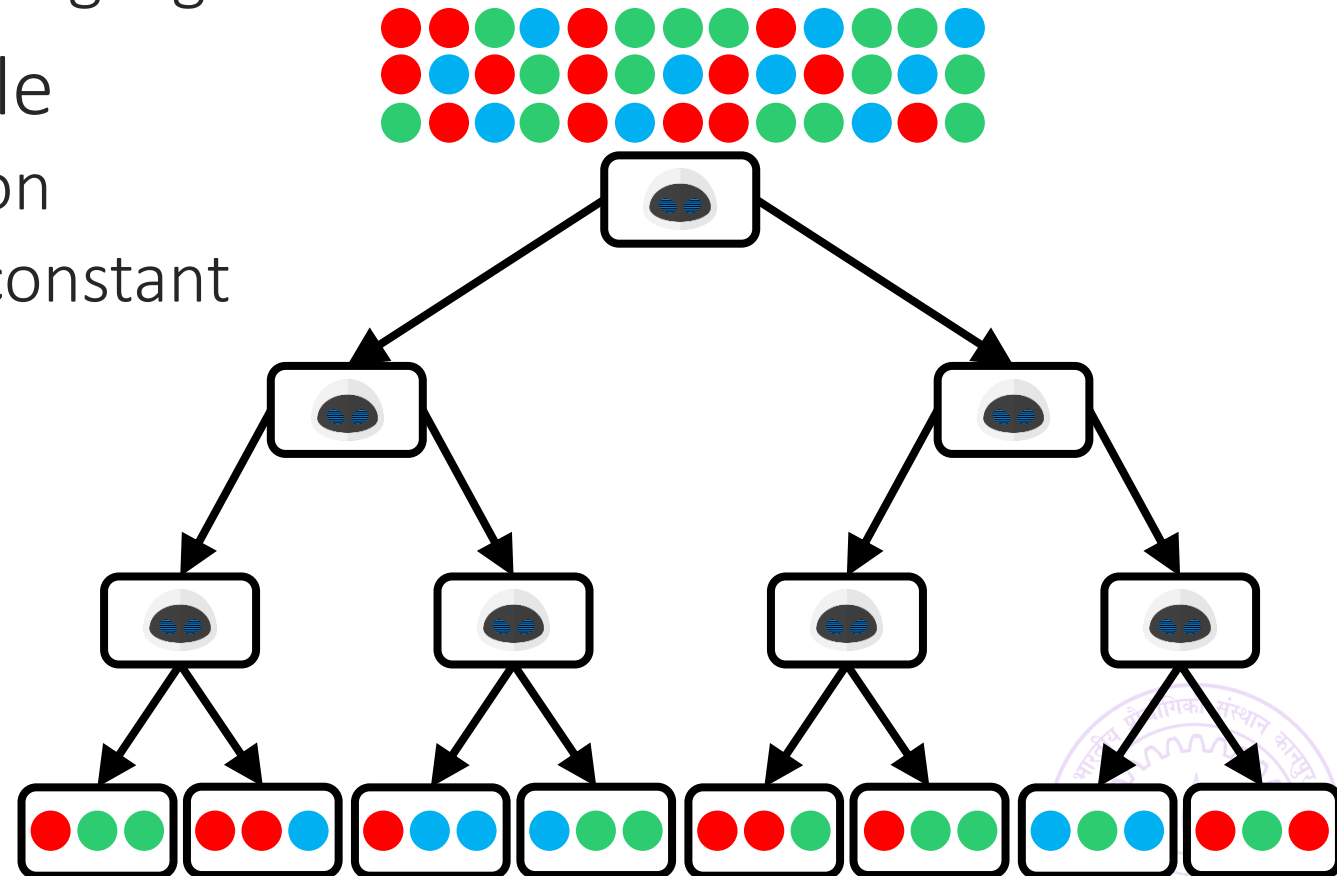
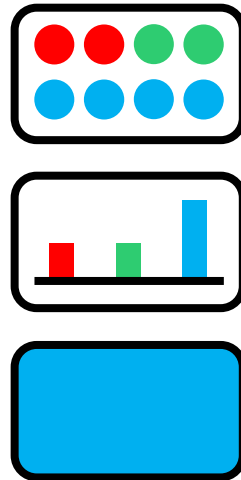
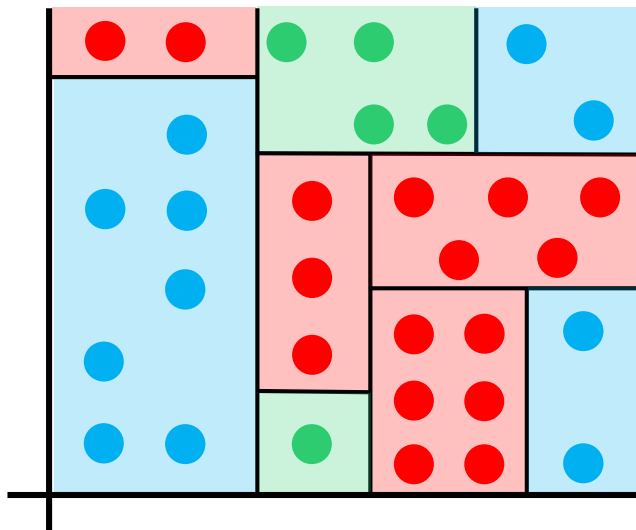
Can take any (complicated) action at a leaf

Why not call another machine learning algorithm?

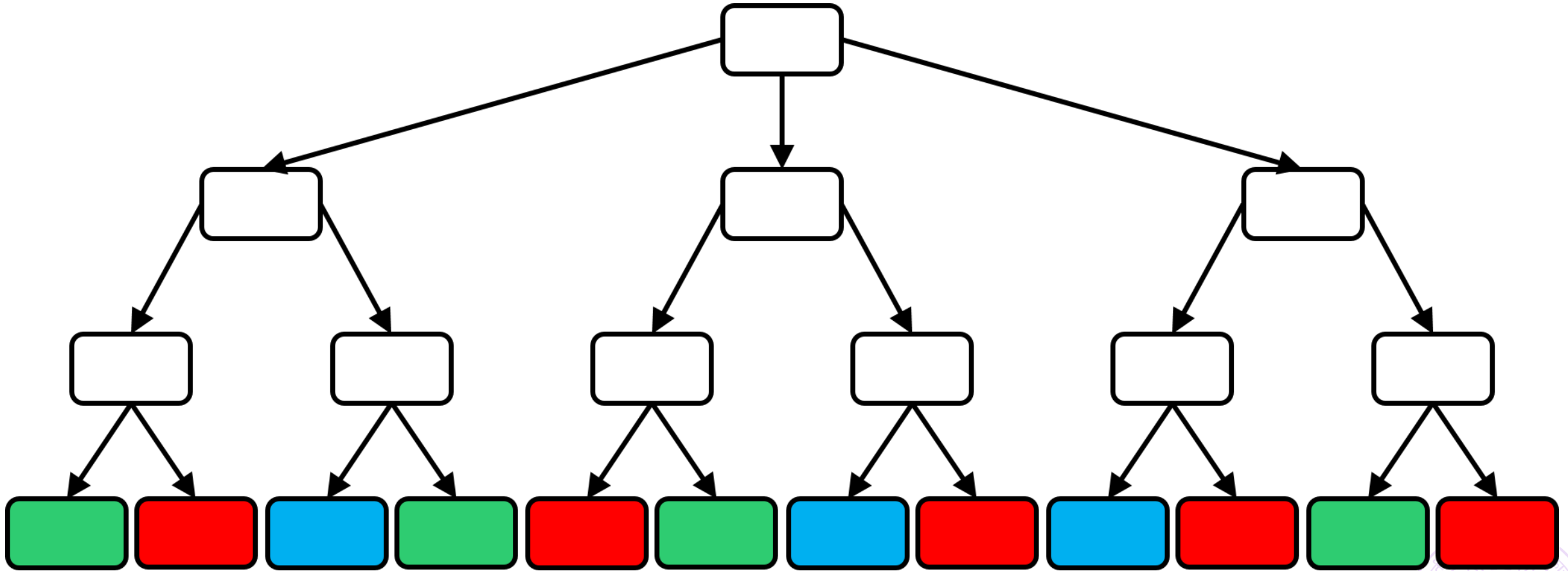
For speed, keep leaf action simple

Simplest action – constant prediction

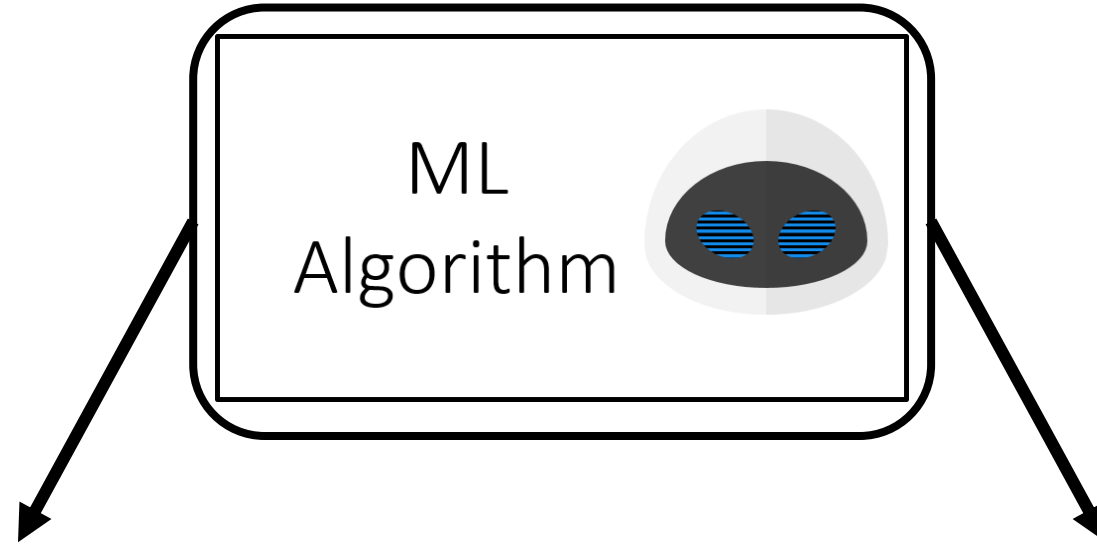
Such a DT will encode a piecewise constant prediction function



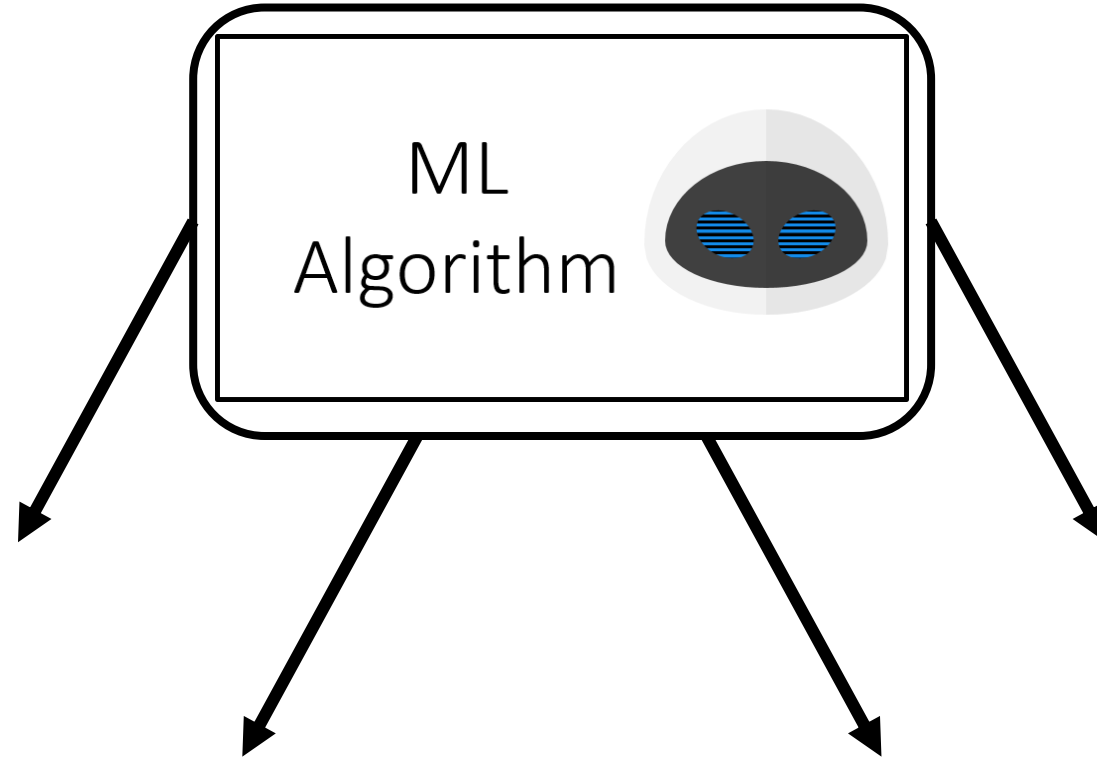
# How to split a node into children nodes? 82



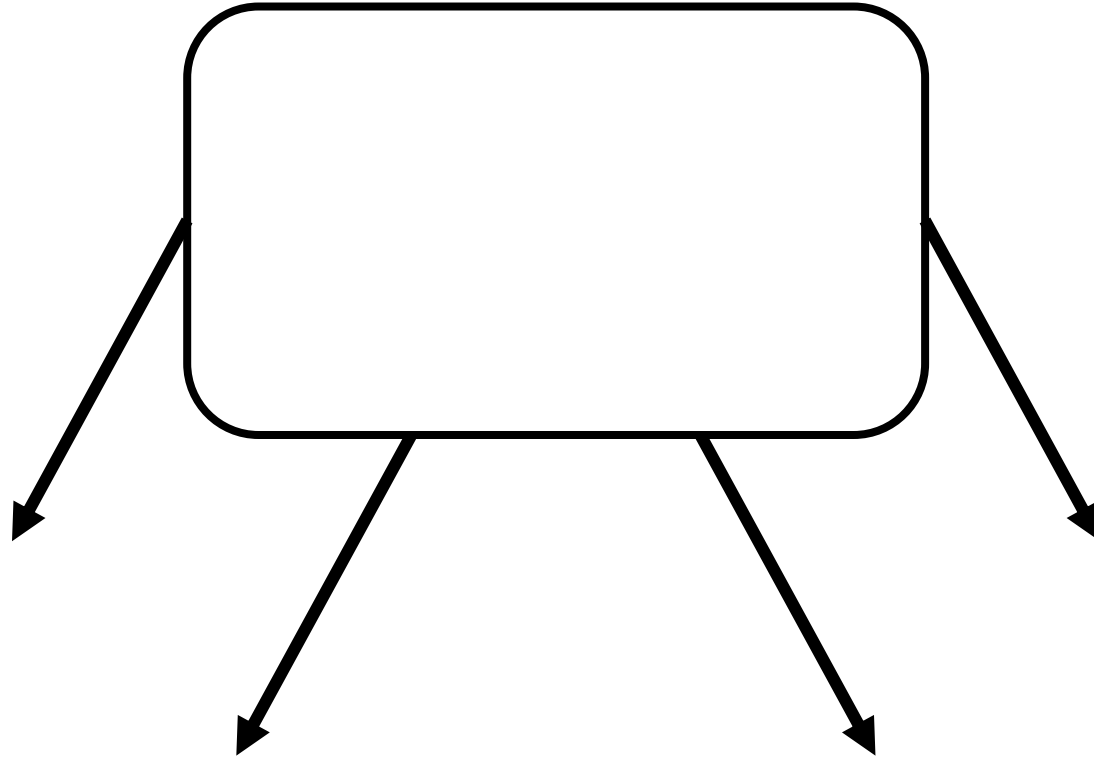
# How to split a node into children nodes? 82



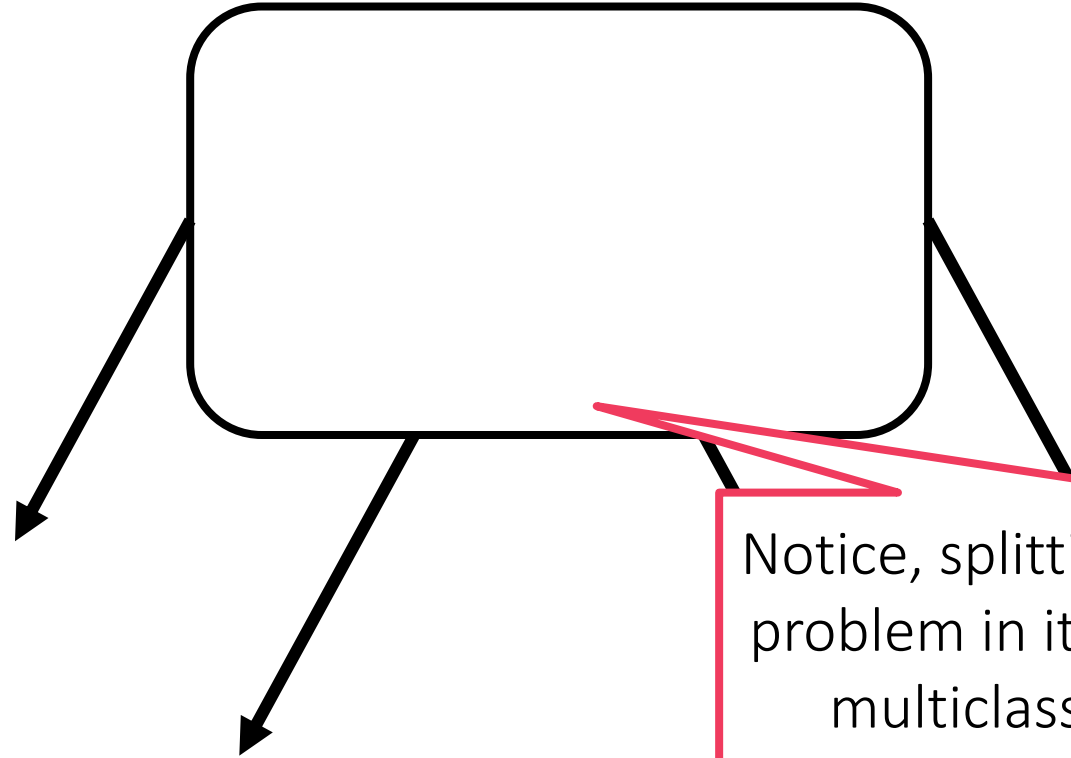
# How to split a node into children nodes? 82



# How to split a node into children nodes? 82



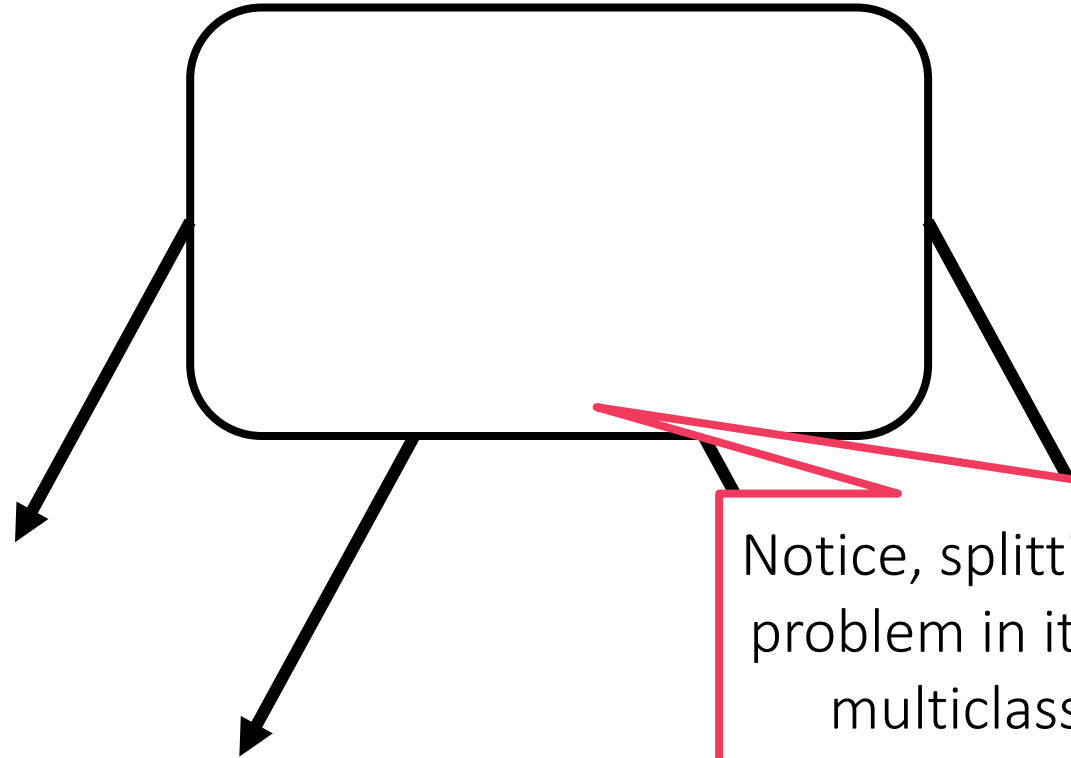
# How to split a node into children nodes? 82



Notice, splitting a node is a classification problem in itself! Binary if two children, multiclass if more than 2 children



# How to split a node into children nodes? 82

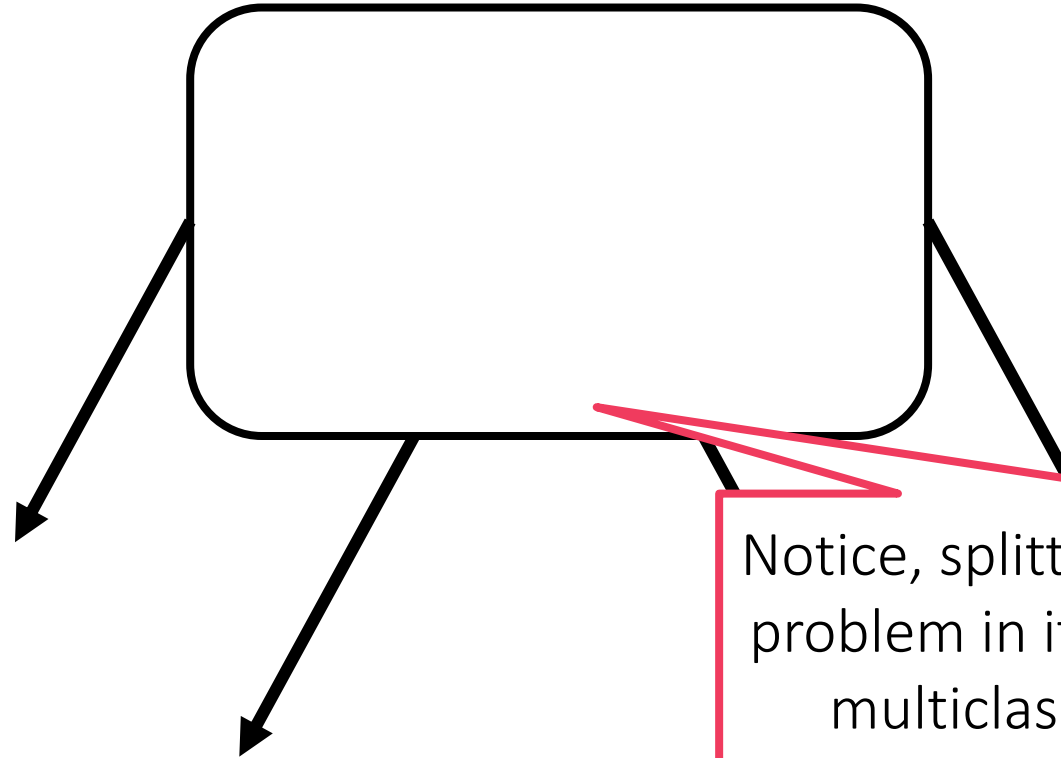


Notice, splitting a node is a classification problem in itself! Binary if two children, multiclass if more than 2 children

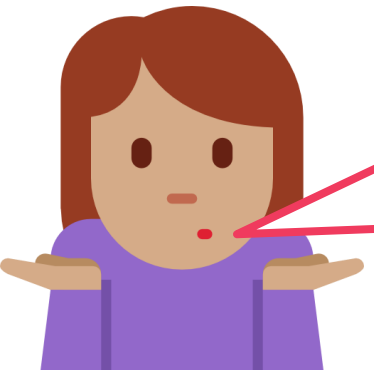
Oh! So we are using a simple ML technique such as binary classification to learn a DT!



# How to split a node into children nodes? 82

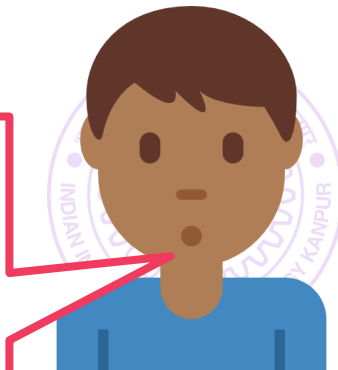


Notice, splitting a node is a classification problem in itself! Binary if two children, multiclass if more than 2 children



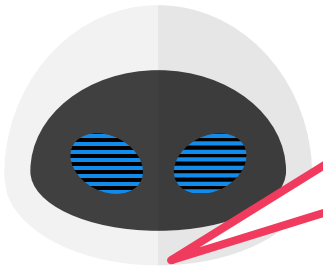
Can we use any classification technique to split a node or are there some restrictions?

Oh! So we are using a simple ML technique such as binary classification to learn a DT!

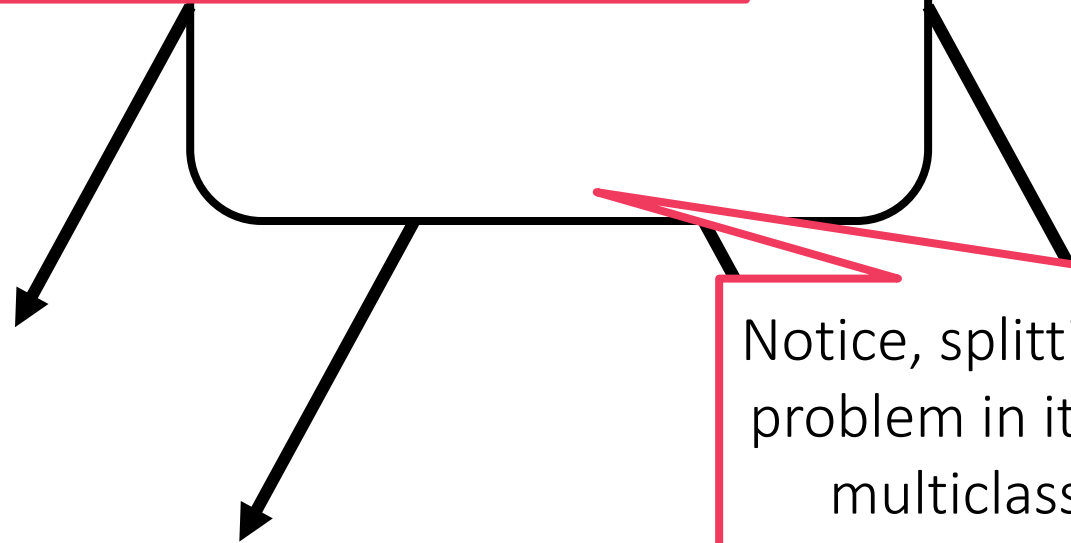




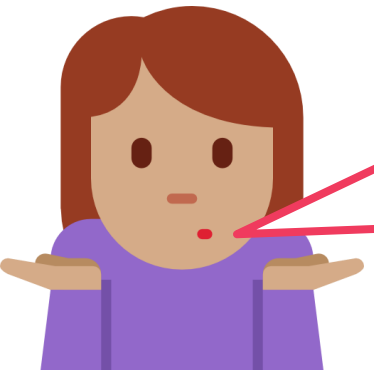
# How to split a node into children nodes? 82



In principle there is no restriction (e.g. can even use a deep net to split a node). However, practical limitations are there.



Notice, splitting a node is a classification problem in itself! Binary if two children, multiclass if more than 2 children



Can we use any classification technique to split a node or are there some restrictions?

Oh! So we are using a simple ML technique such as binary classification to learn a DT!



# Splitting a Node – some lessons

90

Recall, one of the goals of DT is to speed up kNN prediction time

Thus, node splitting algorithm must be super fast otherwise no benefit of using DT – may just as well do NN directly

Often people carefully choose just a single feature and split a node based on that e.g. (age  $< 25$  go left, age  $\geq 25$  go right)

Such “simple classifiers” are often called *decision stumps*



# Splitting a Node – s

Recall, one of the goals of DT

Thus, the splitting algorithm must be super fast otherwise no

ben  
Often, making sure that the split is balanced (e.g. roughly half the data points go left and right) is also important to ensure that the tree is balanced. However, Such ensuring balance is often tricky

Various notions of purity exist – entropy and Gini index for classification problems, variance for regression problems

Pure nodes are very convenient. We can make them leaves right away and not have to worry about splitting them 😊

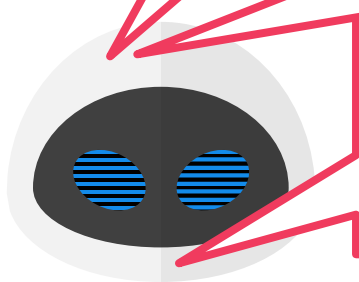
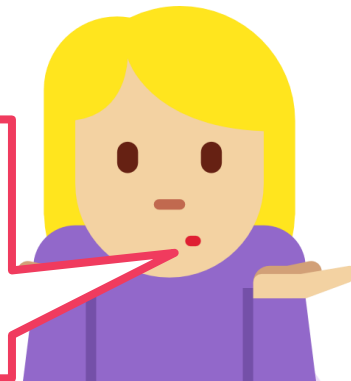
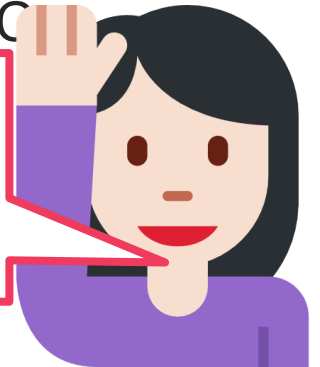
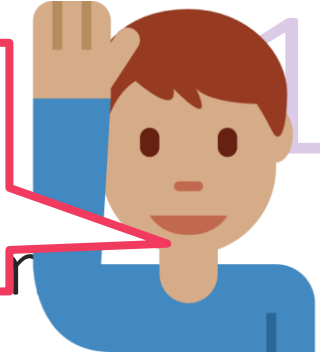
go left, age  $\leq 25$  go right)

often called *decision stumps*

A child node is completely pure if it contains training data of only one class.

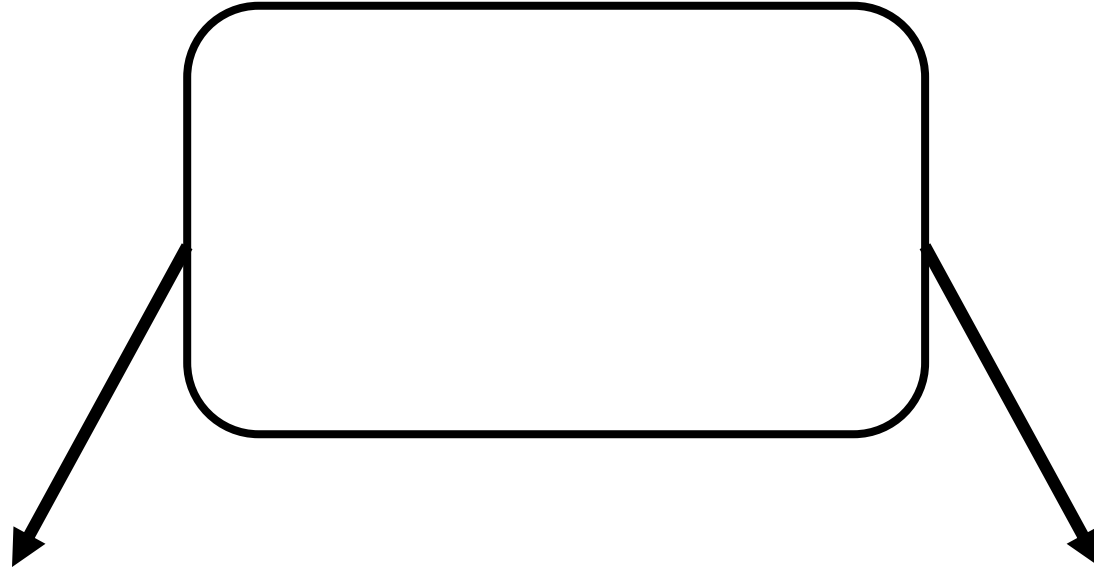
How do I decide whether to use age or gender? Even if using age, how do I decide whether to threshold at 25 or 65?

Usually, people go over all available features and all possible thresholds (can be slow if not done cleverly) and choose a feature and a threshold for that feature so that the child nodes that are created are as *pure* as possible



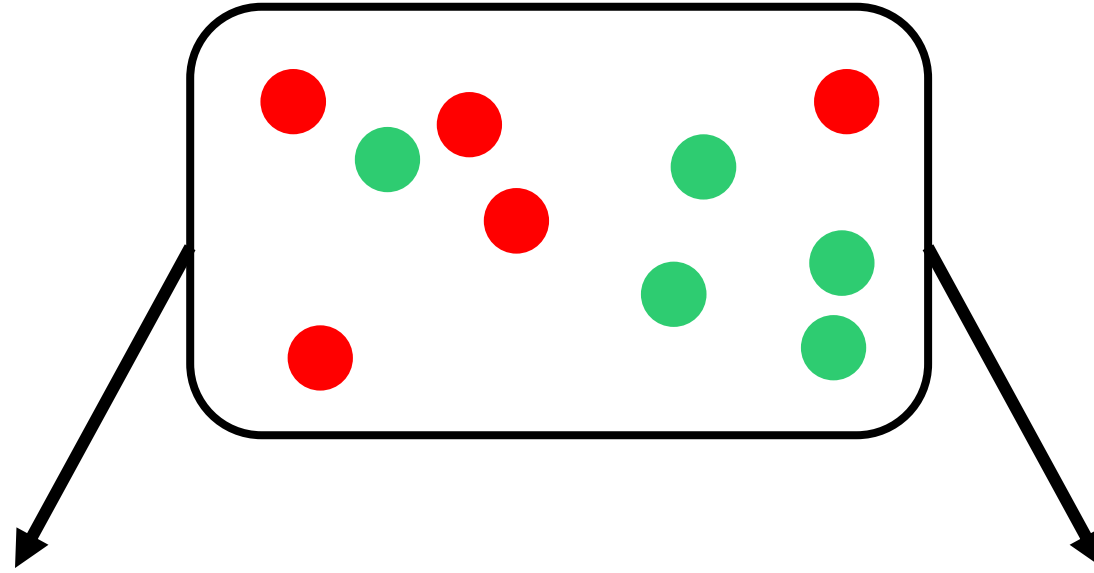
# Purifying Decision Stumps

92



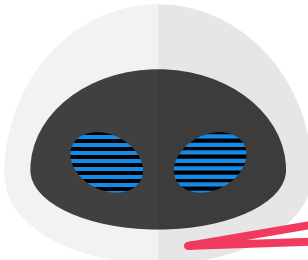
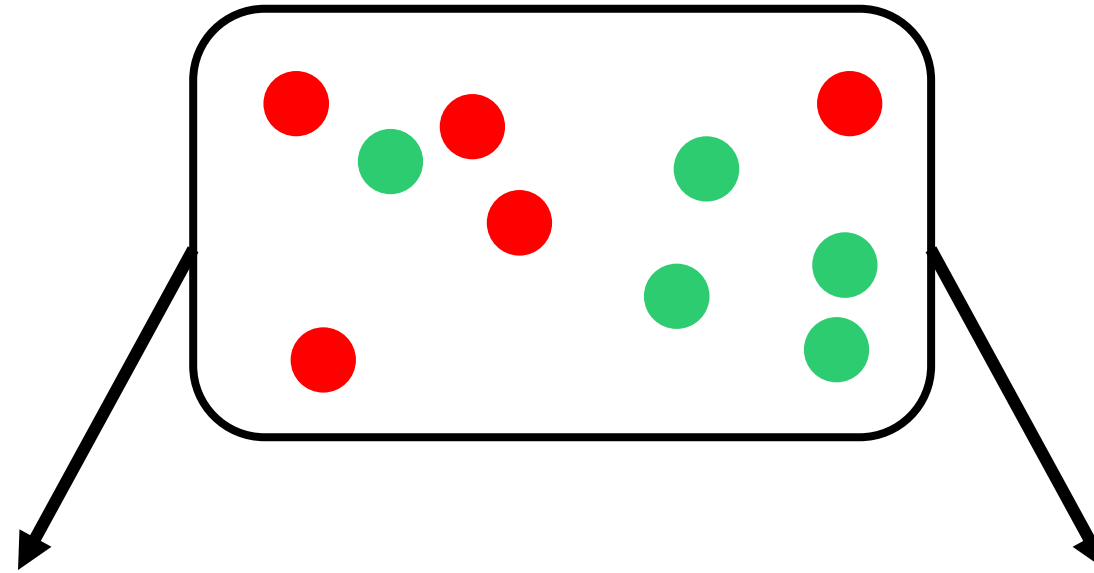
# Purifying Decision Stumps

92



# Purifying Decision Stumps

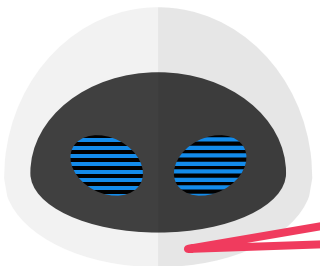
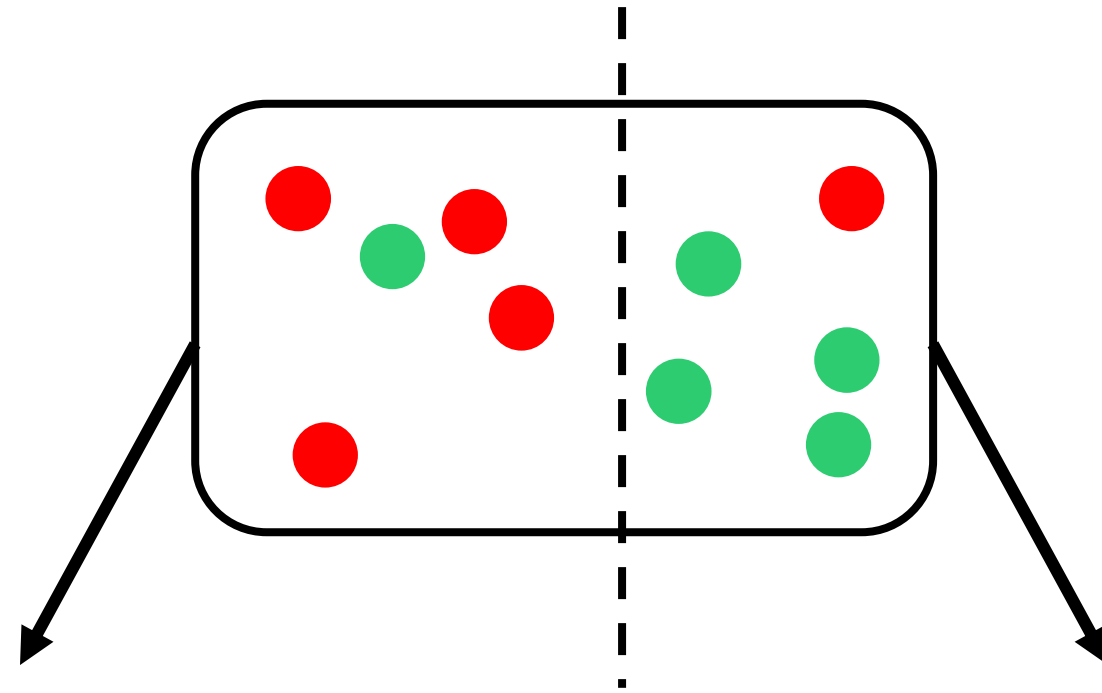
92



Two possible splitting directions. Let us choose the one that gives us purer children

# Purifying Decision Stumps

92

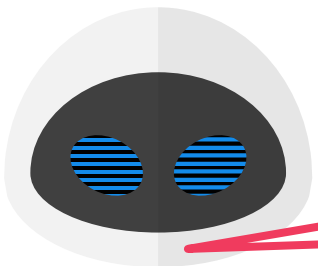
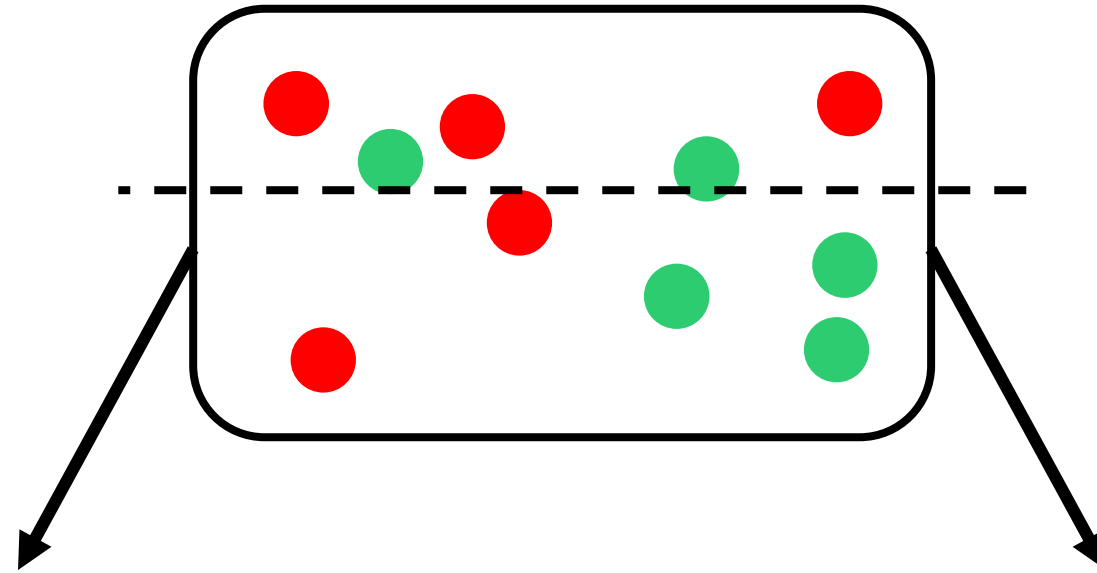


Two possible splitting directions. Let us choose the one that gives us purer children



# Purifying Decision Stumps

92



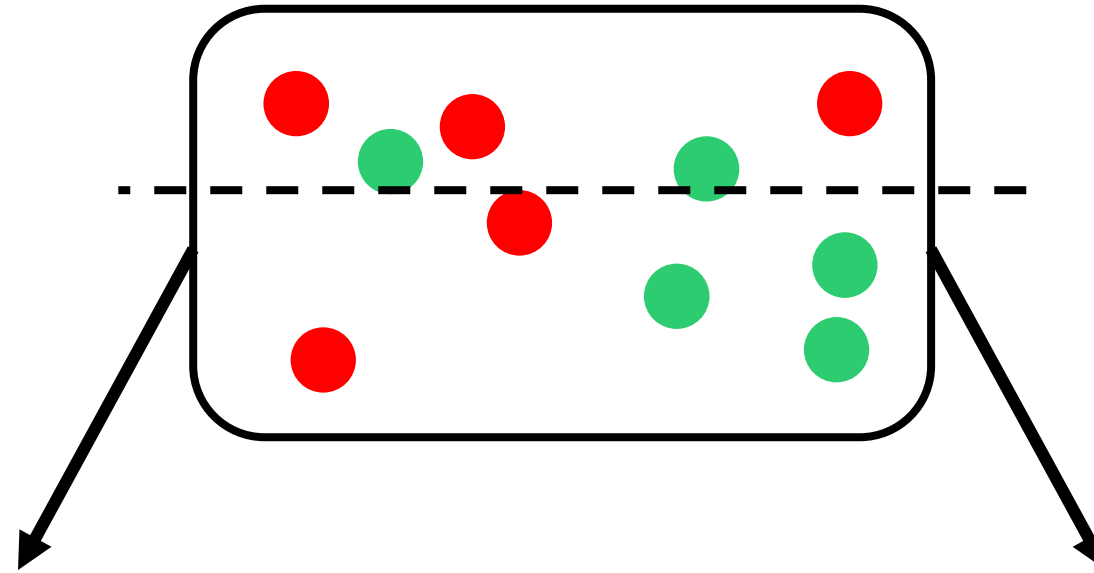
Two possible splitting directions. Let us choose the one that gives us purer children





# Purifying Decision Stumps

92



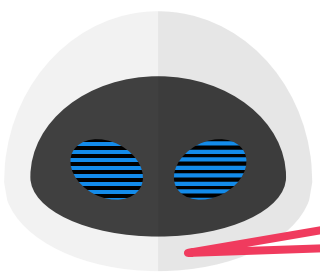
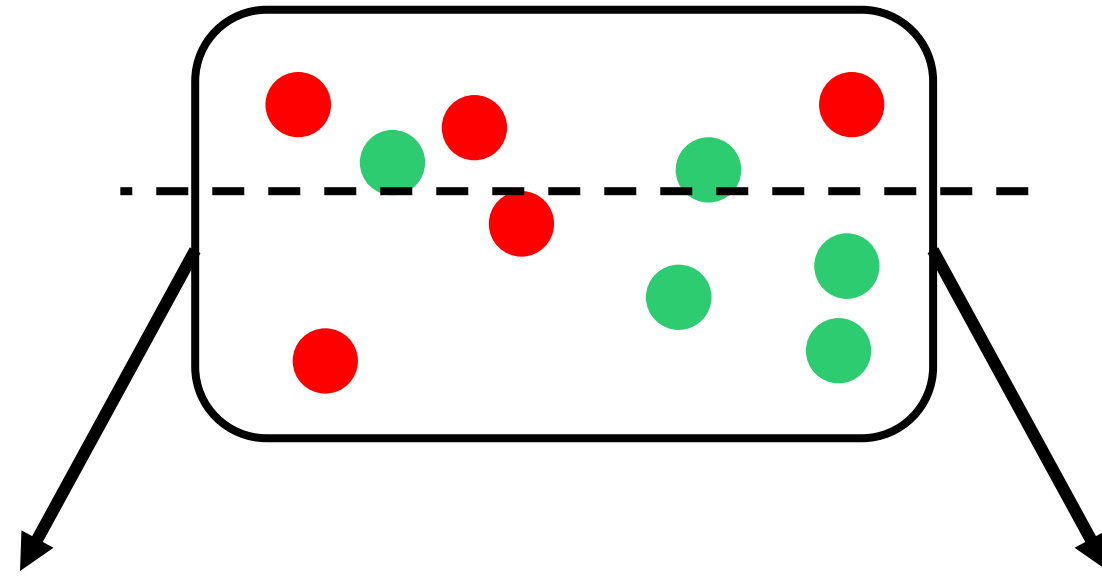
Search purest horizontal split by going over all possible thresholds and checking

Two possible splitting directions. Let us choose the one that gives us purer children



# Purifying Decision Stumps

92

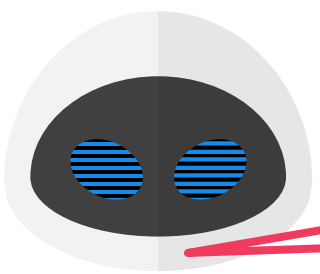
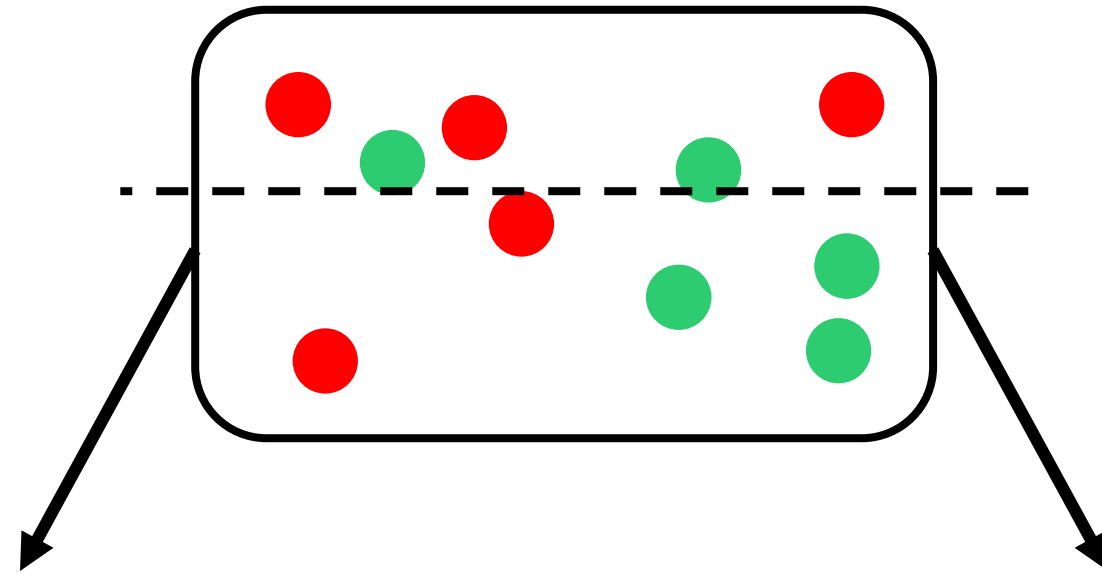


Two possible splitting directions. Let us choose the one that gives us purer children



# Purifying Decision Stumps

92



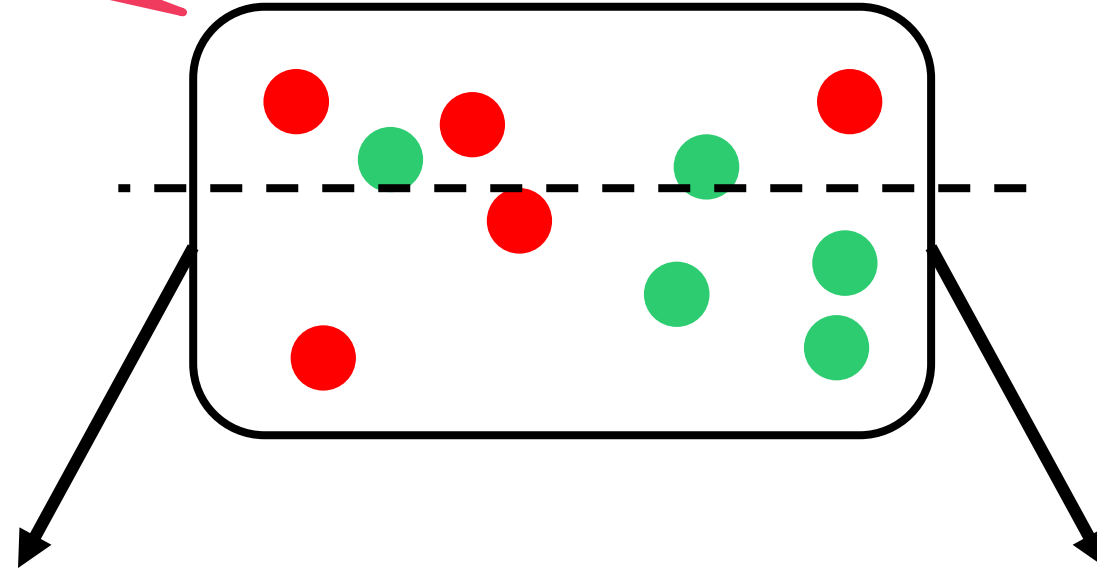
Two possible splitting directions. Let us choose the one that gives us purer children



# Purifying Decision Stumps

92

Purest horizontal split



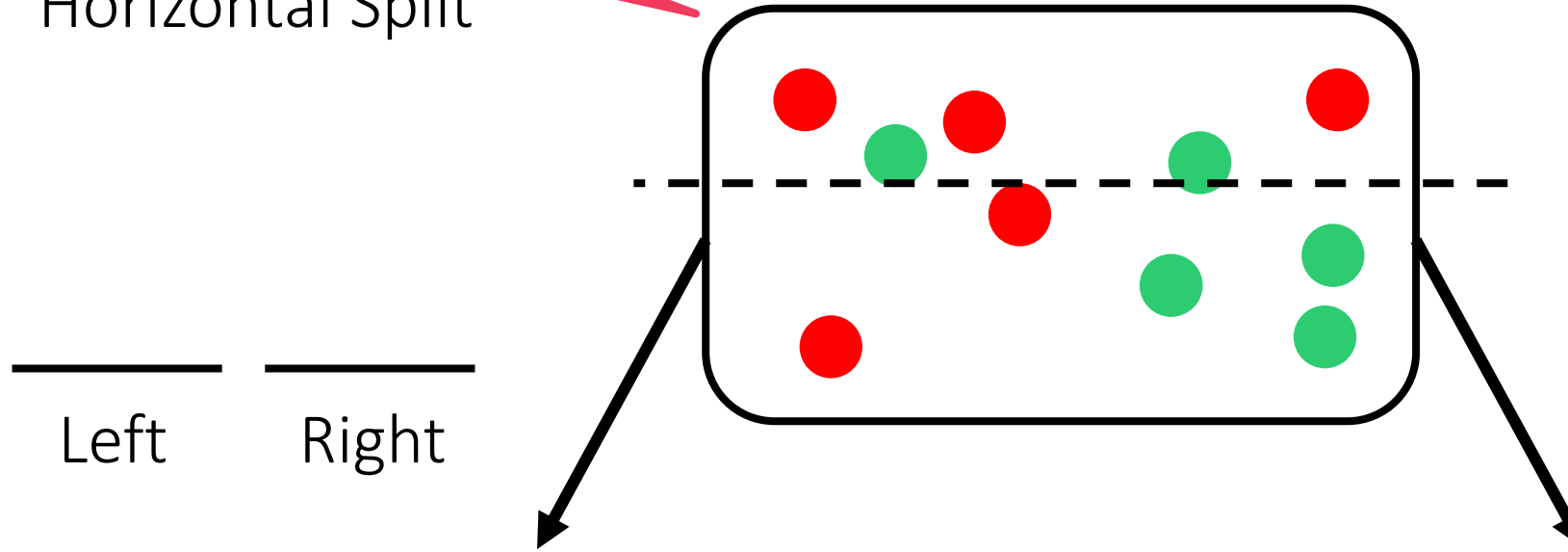
Two possible splitting directions. Let us choose the one that gives us purer children



# Purifying Decision Stumps

Purest horizontal split

Horizontal Split

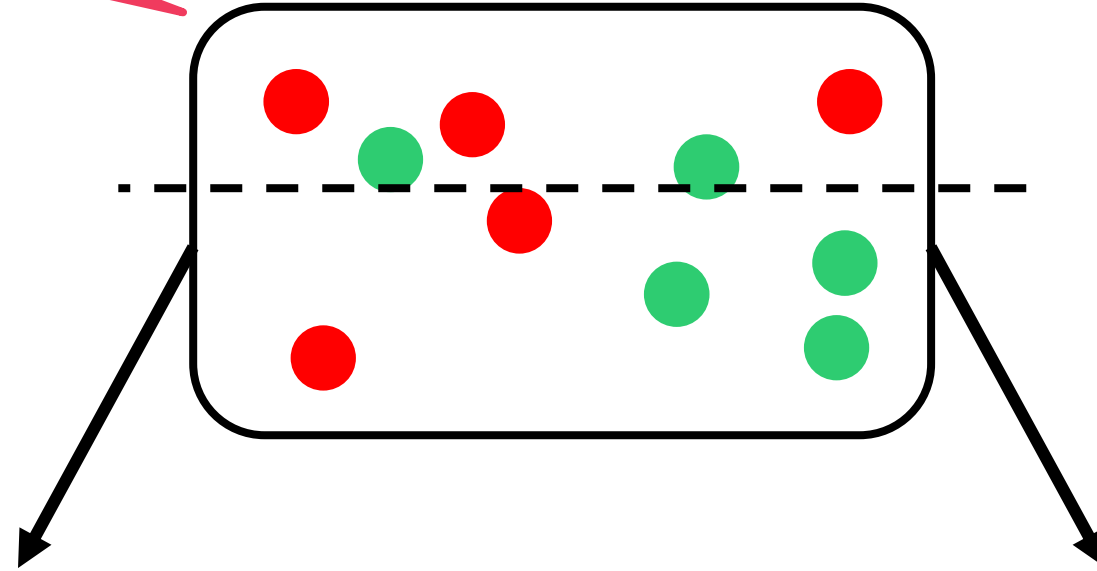
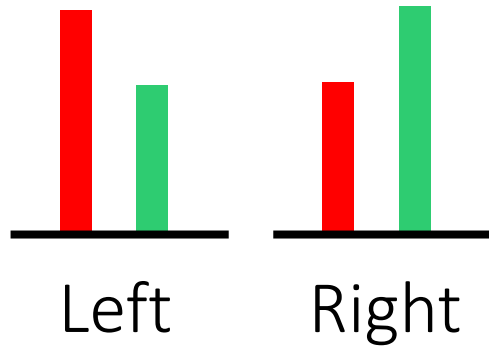


Two possible splitting directions. Let us choose the one that gives us purer children

# Purifying Decision Stumps

Purest horizontal split

Horizontal Split

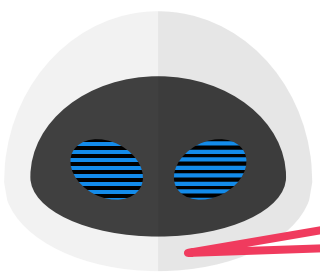
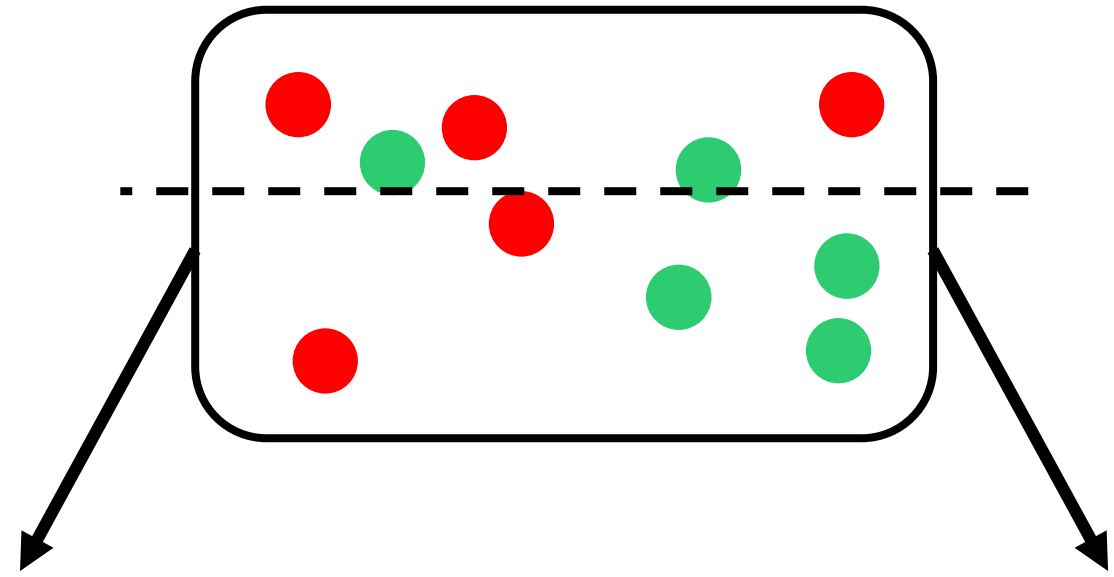
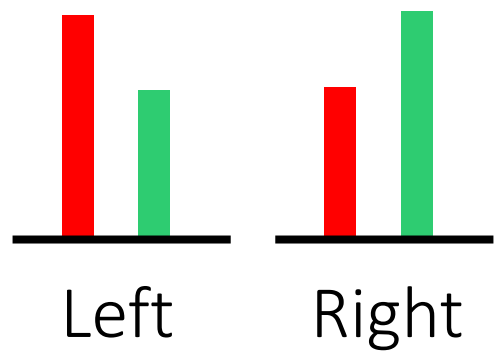


Two possible splitting directions. Let us choose the one that gives us purer children

# Purifying Decision Stumps

92

Horizontal Split



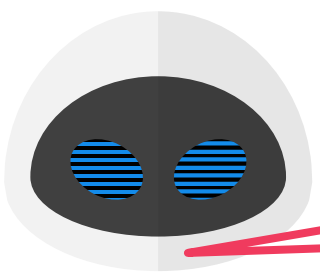
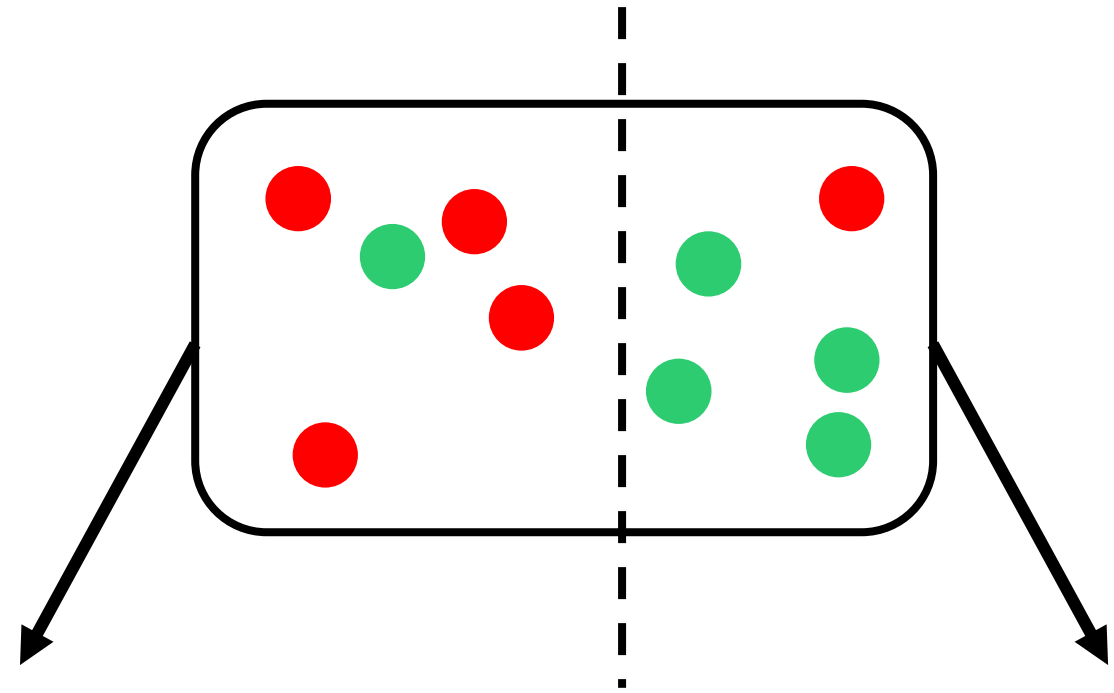
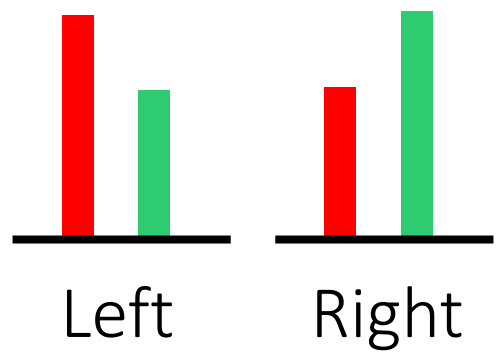
Two possible splitting directions. Let us choose the one that gives us purer children



# Purifying Decision Stumps

92

Horizontal Split



Two possible splitting directions. Let us choose the one that gives us purer children

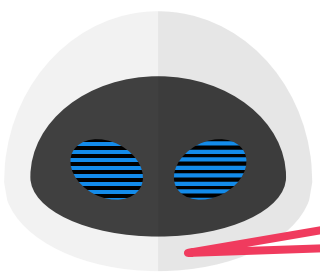
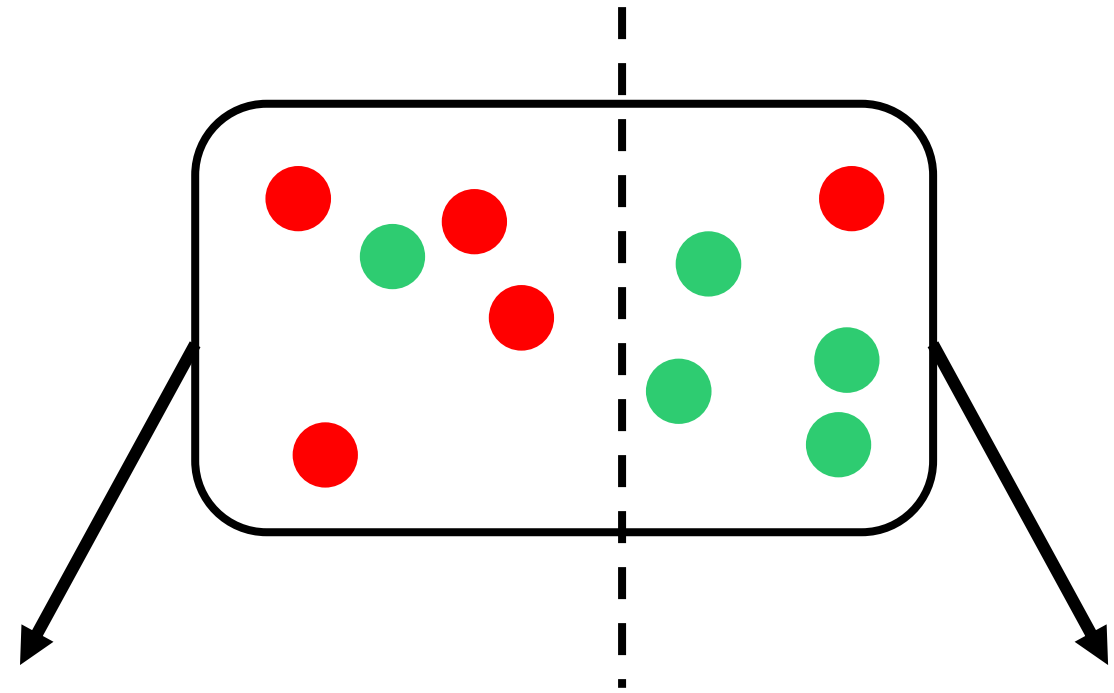
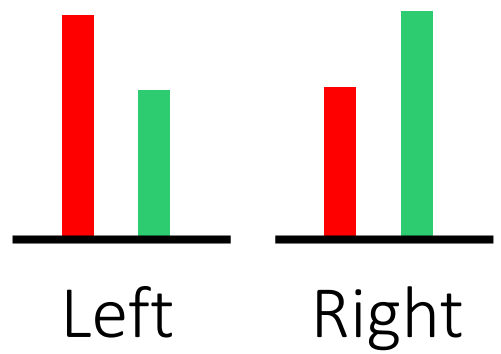




# Purifying Decision Stumps

92

Horizontal Split

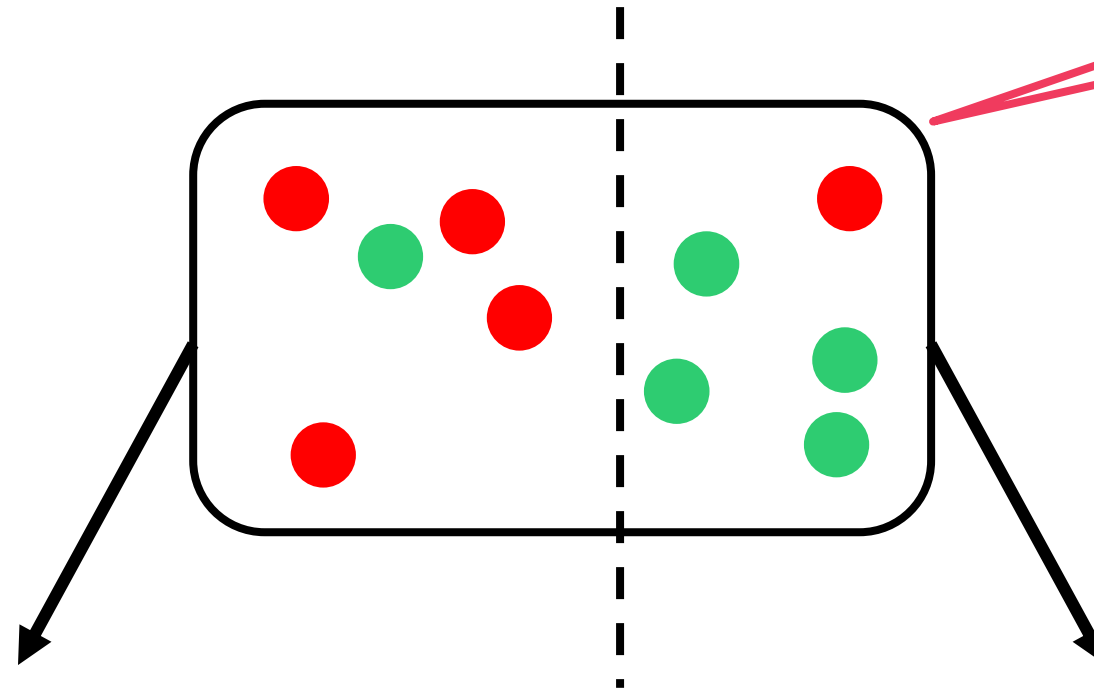
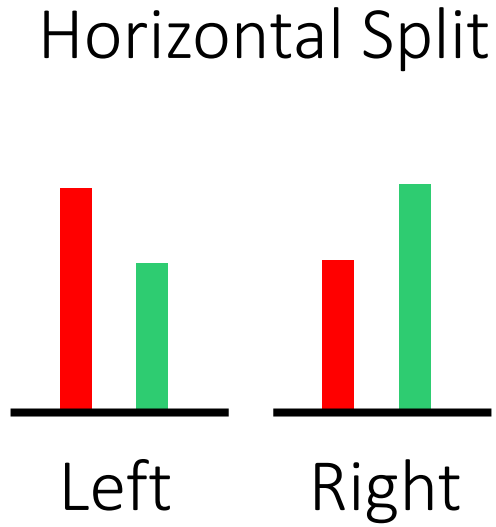


Two possible splitting directions. Let us choose the one that gives us purer children

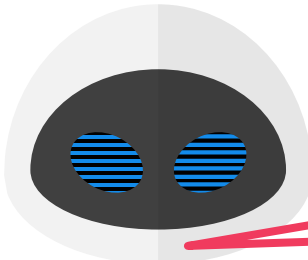


# Purifying Decision Stumps

02



Purest vertical split



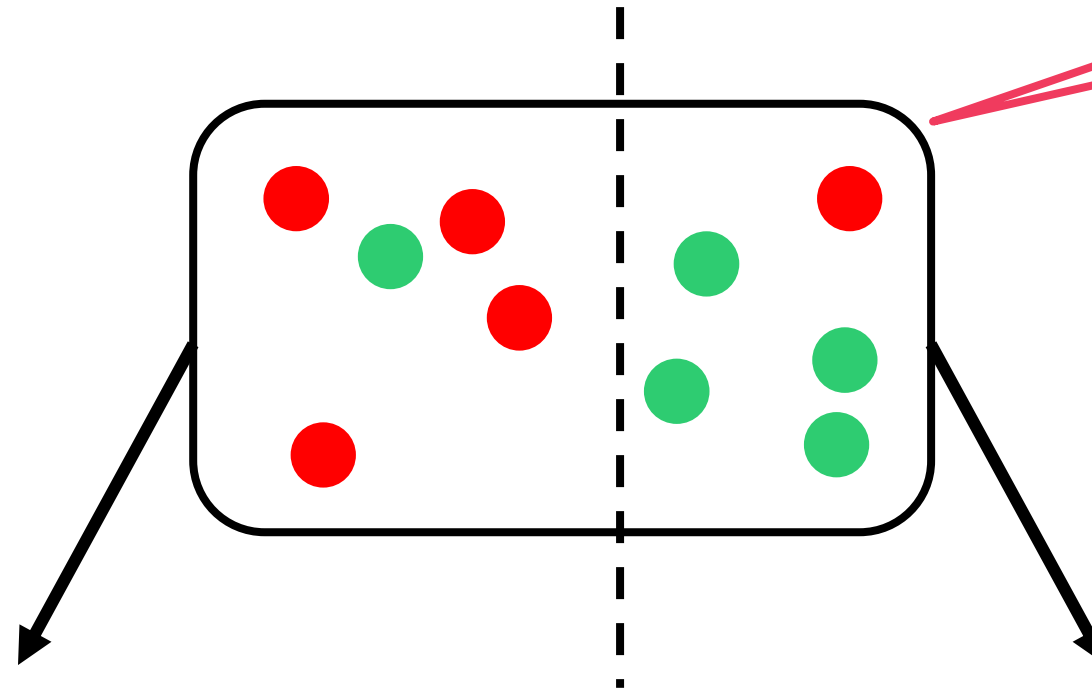
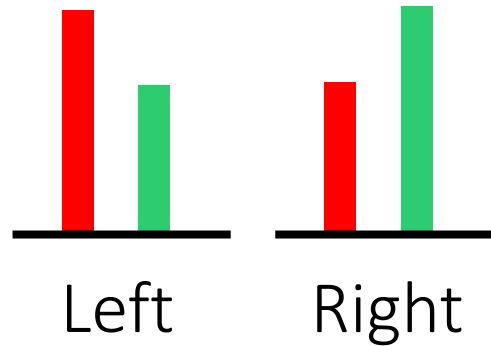
Two possible splitting directions. Let us choose the one that gives us purer children



# Purifying Decision Stumps

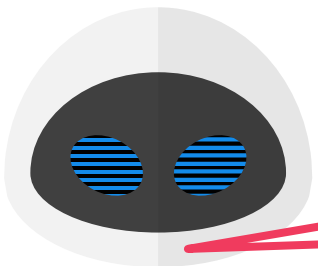
02

Horizontal Split



Purest vertical split

Vertical Split

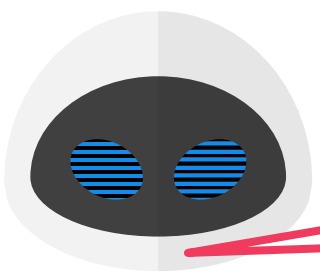
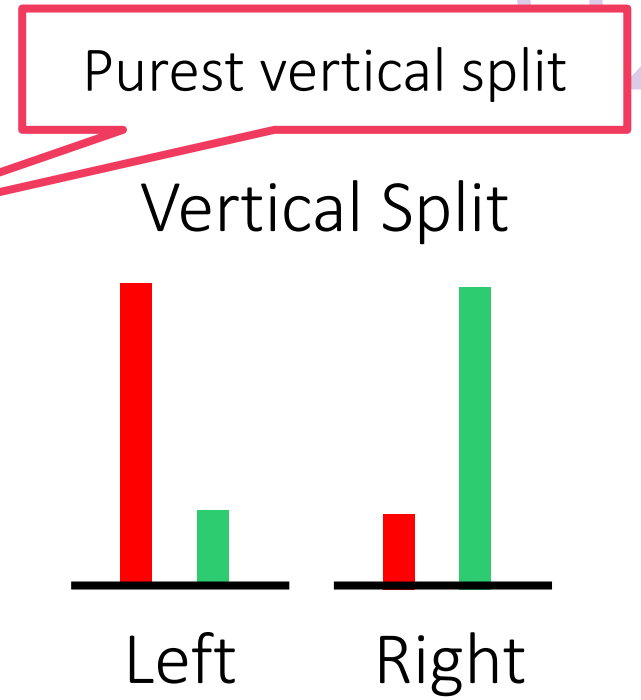
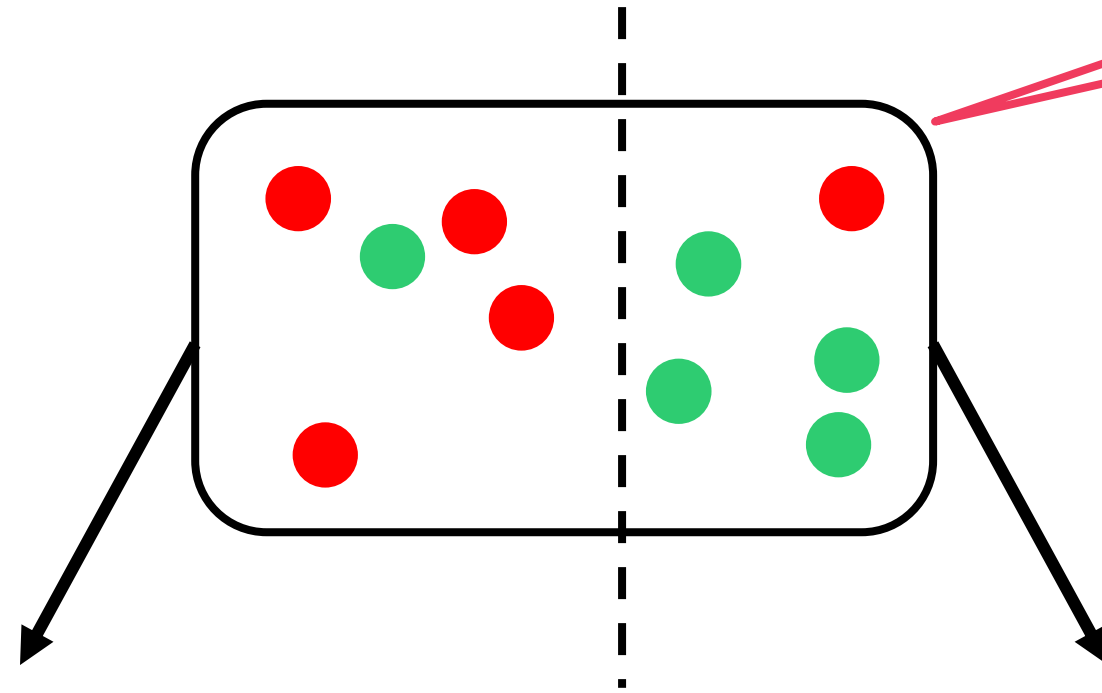
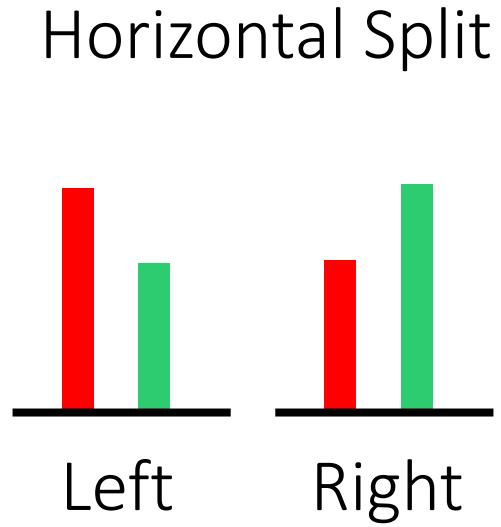


Two possible splitting directions. Let us choose the one that gives us purer children



# Purifying Decision Stumps

02



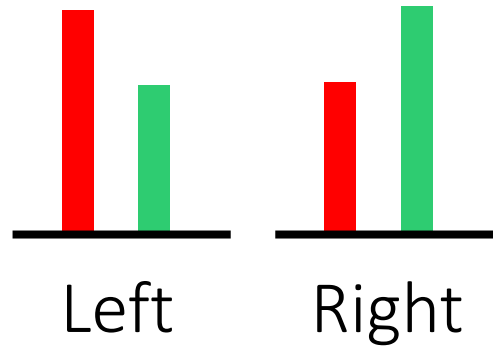
Two possible splitting directions. Let us choose the one that gives us purer children



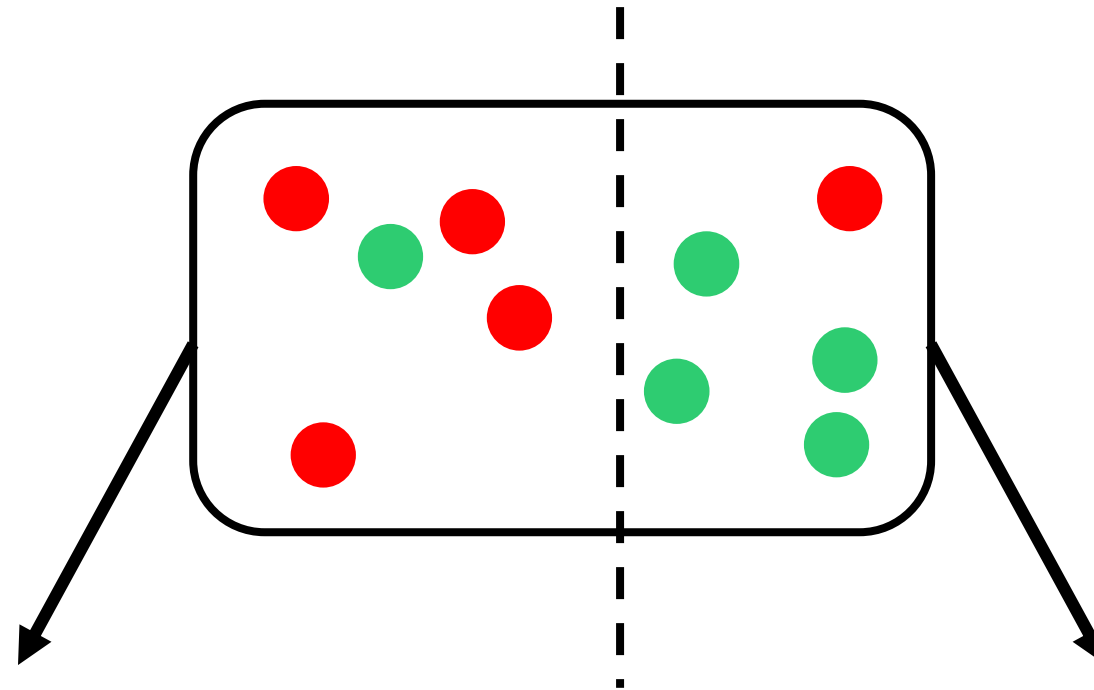
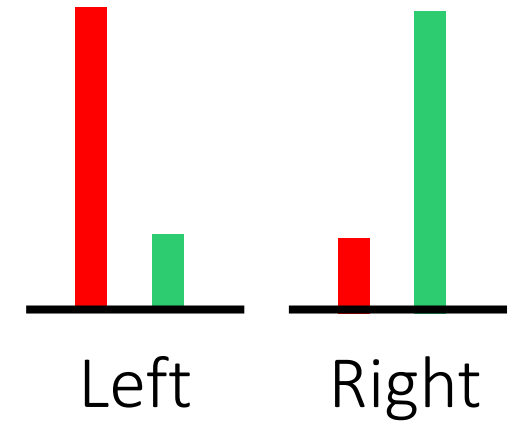
# Purifying Decision Stumps

92

Horizontal Split



Vertical Split

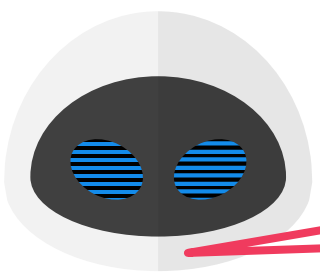
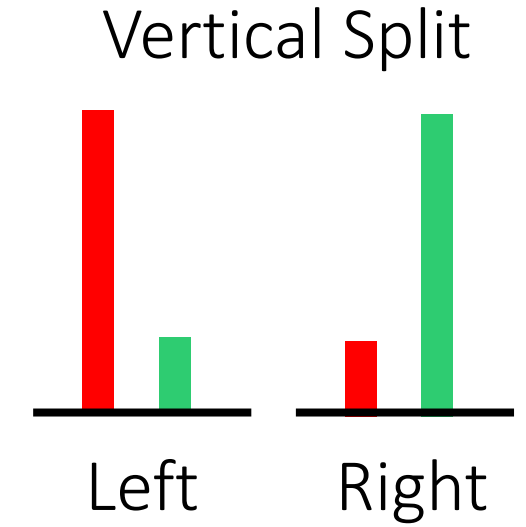
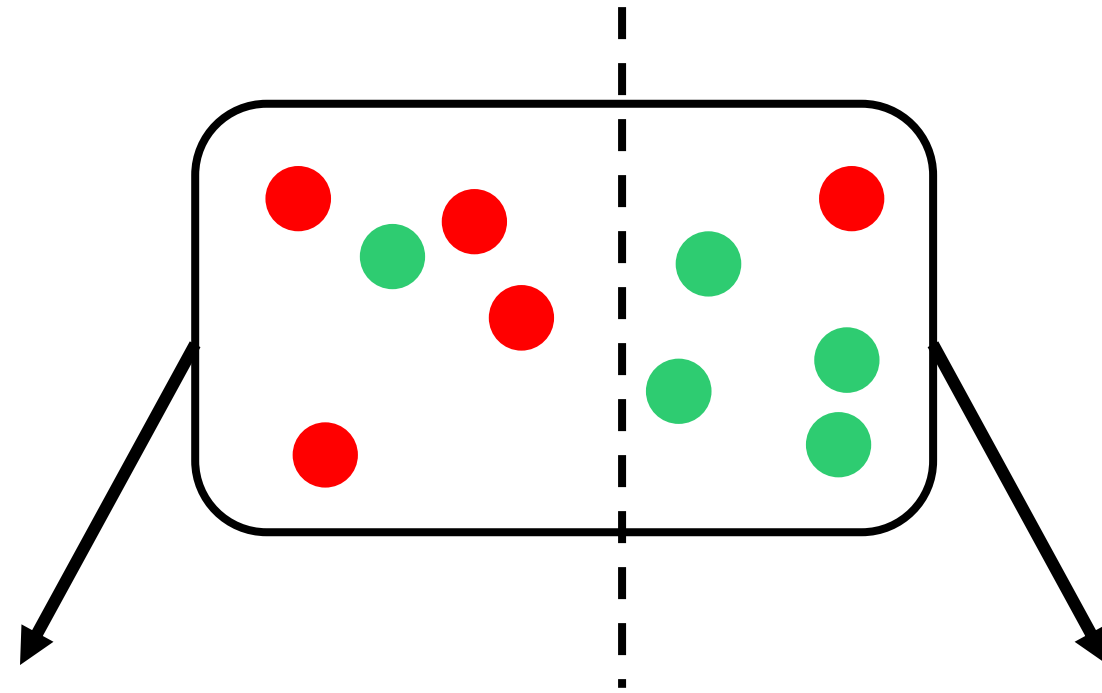
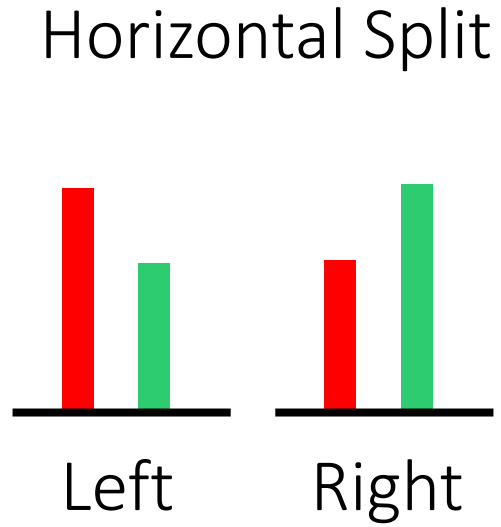


Two possible splitting directions. Let us choose the one that gives us purer children



# Purifying Decision Stumps

92



Two possible splitting directions. Let us choose the one that gives us purer children

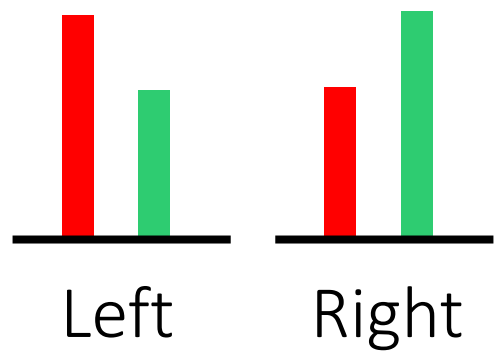
Vertical split is more pure – lets go with it!



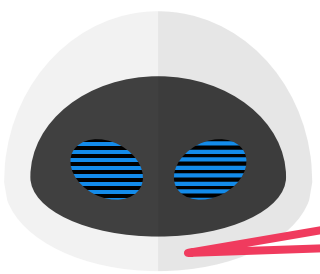
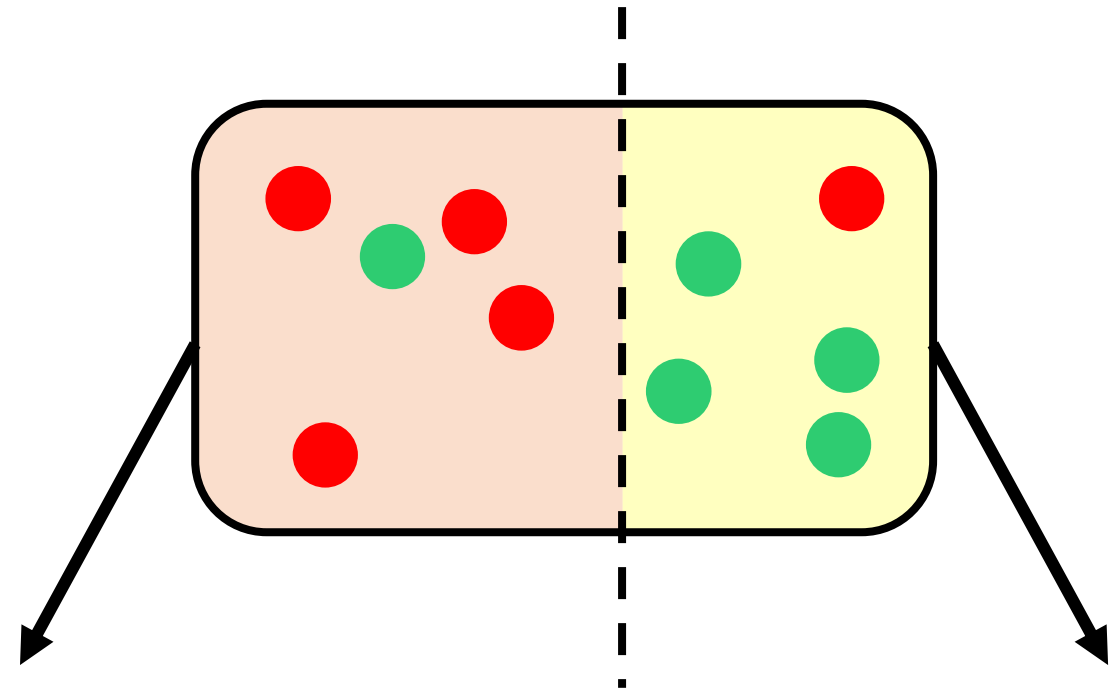
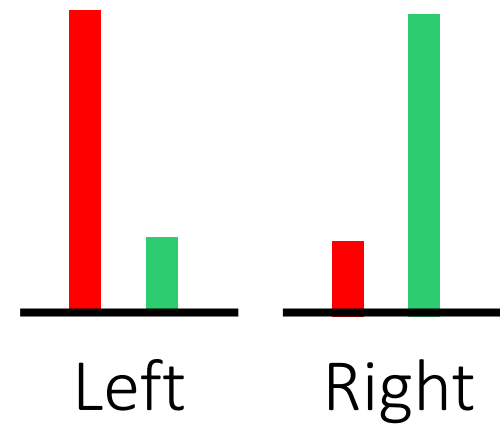
# Purifying Decision Stumps

92

Horizontal Split



Vertical Split



Two possible splitting directions. Let us choose the one that gives us purer children

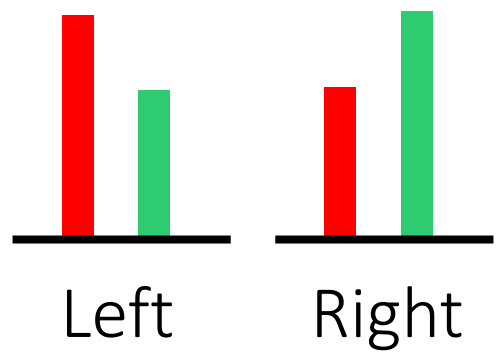
Vertical split is more pure – lets go with it!



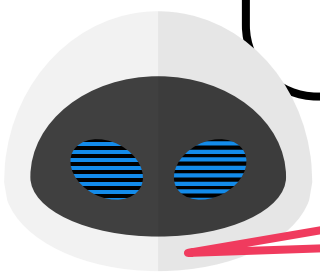
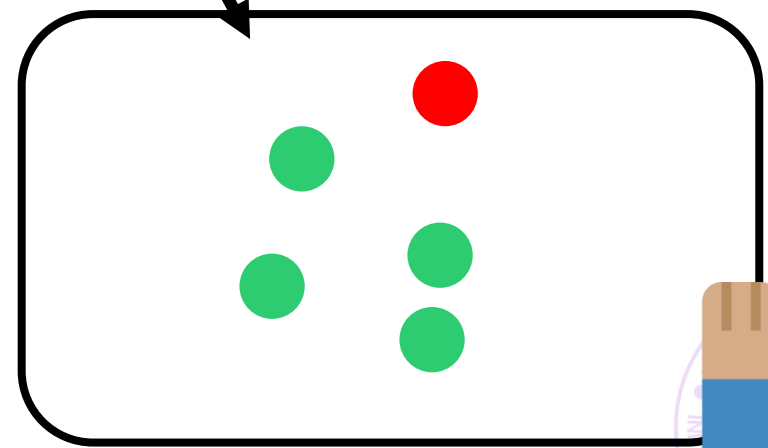
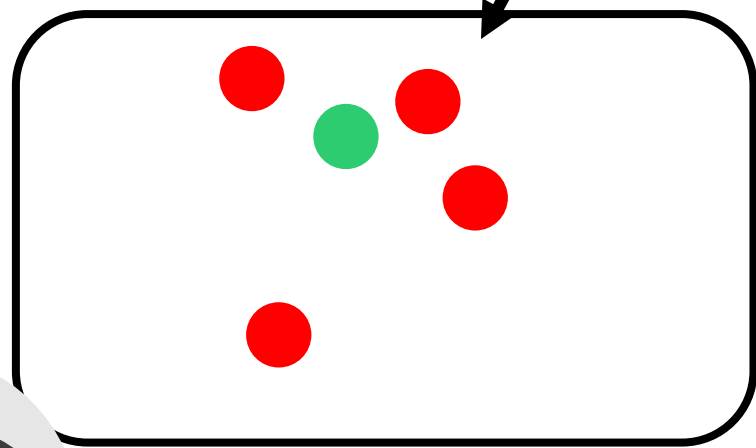
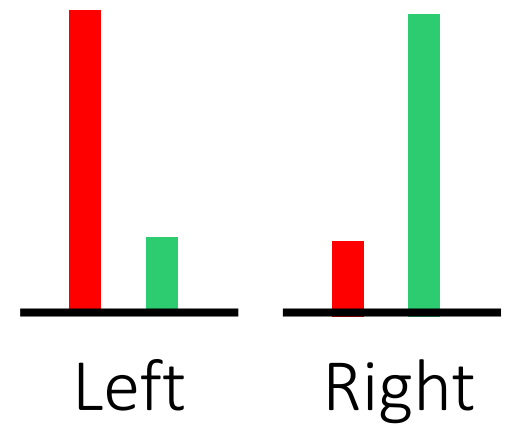
# Purifying Decision Stumps

92

Horizontal Split



Vertical Split



Two possible splitting directions. Let us choose the one that gives us purer children

Vertical split is more pure – lets go with it!

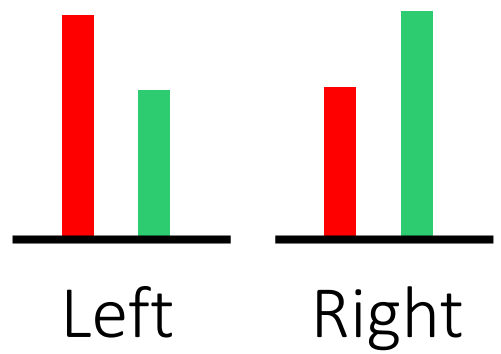




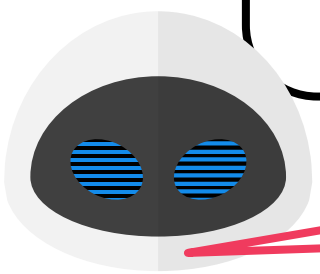
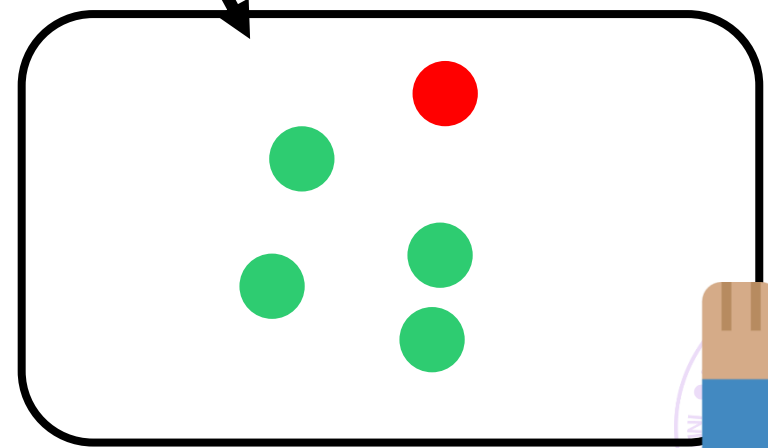
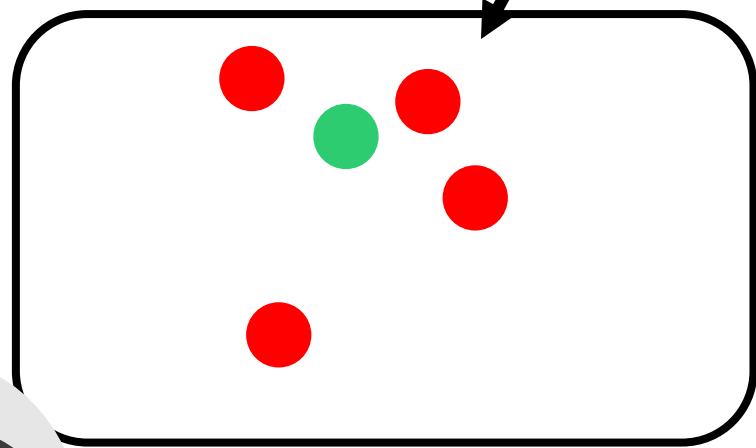
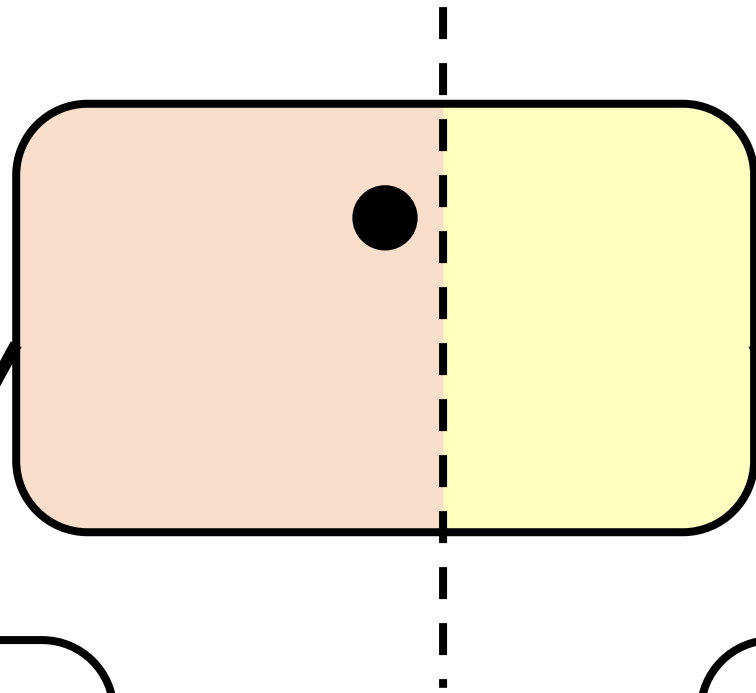
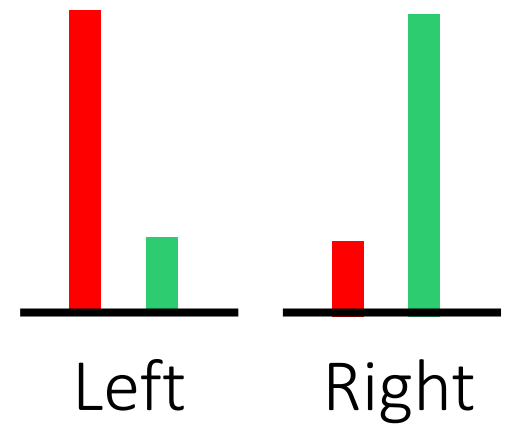
# Purifying Decision Stumps

92

Horizontal Split



Vertical Split



Two possible splitting directions. Let us choose the one that gives us purer children

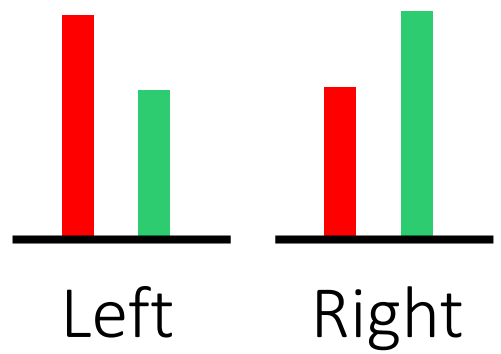
Vertical split is more pure – lets go with it!



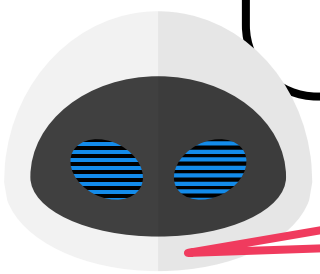
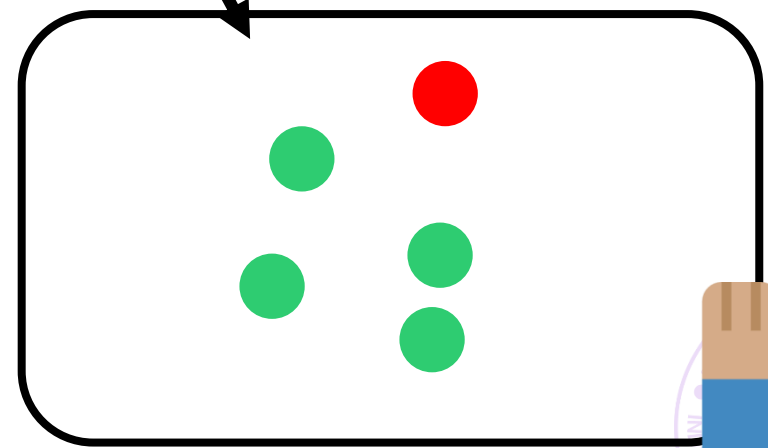
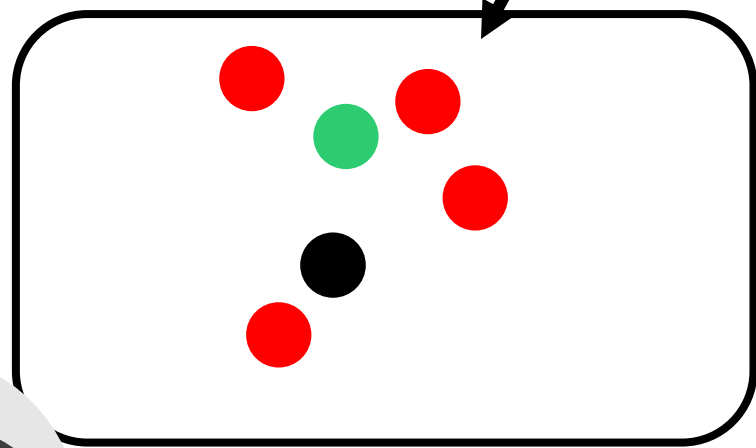
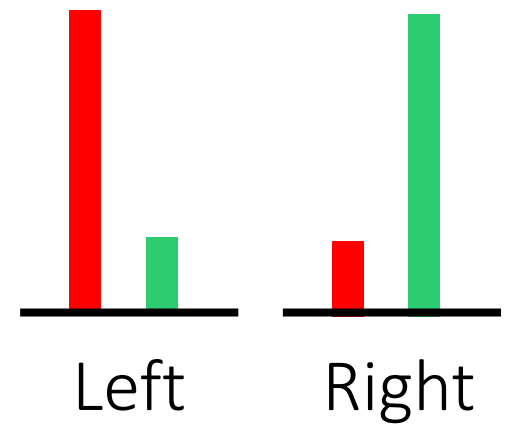
# Purifying Decision Stumps

92

Horizontal Split



Vertical Split



Two possible splitting directions. Let us choose the one that gives us purer children

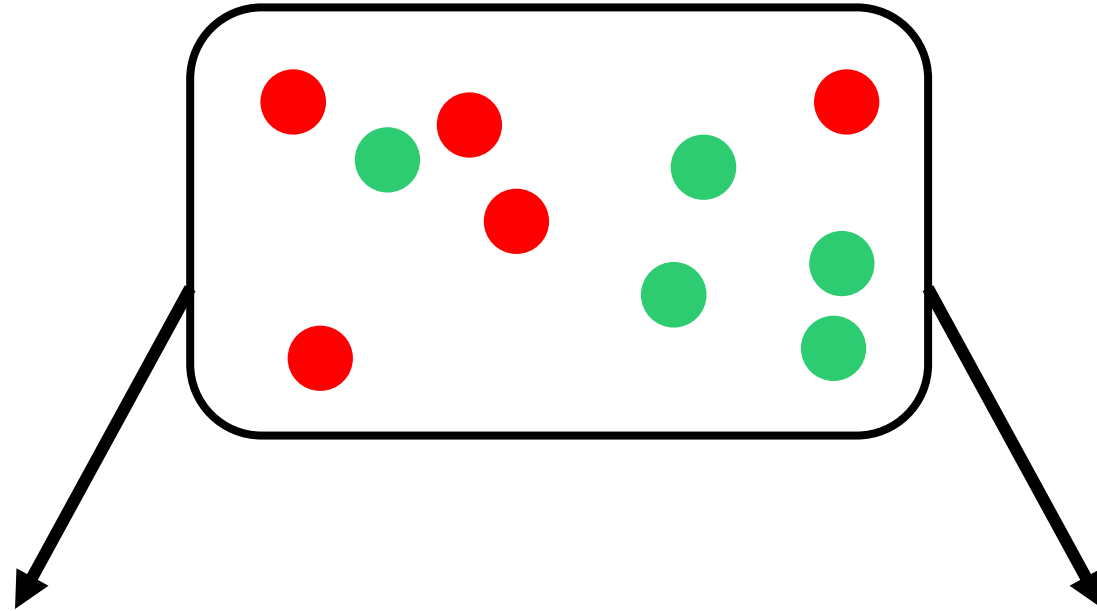
Vertical split is more pure – lets go with it!



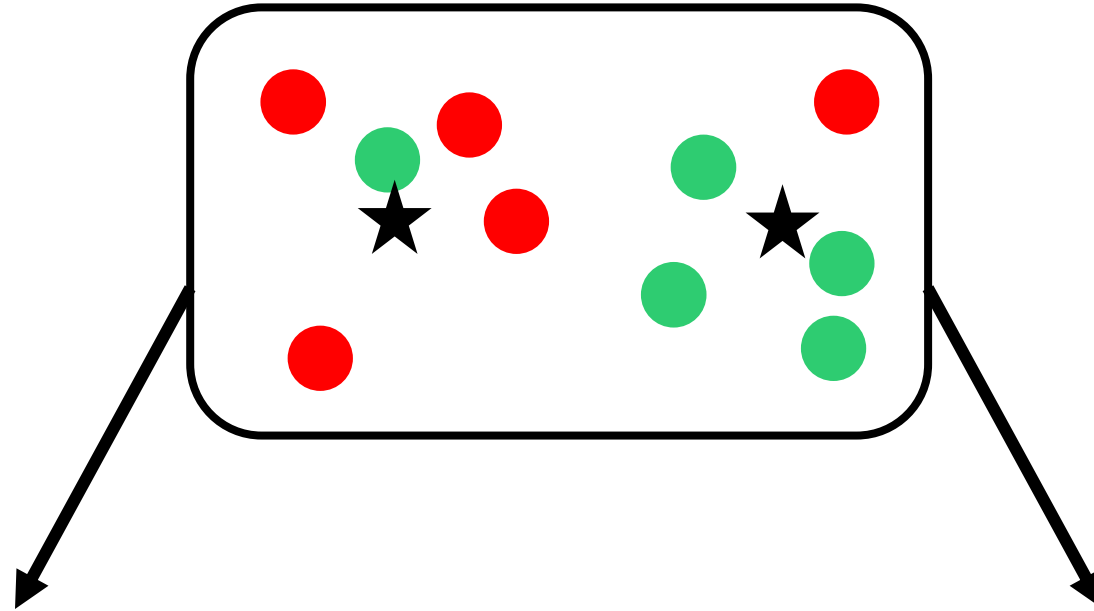
# Node splitting via LwP/linear classifiers 115



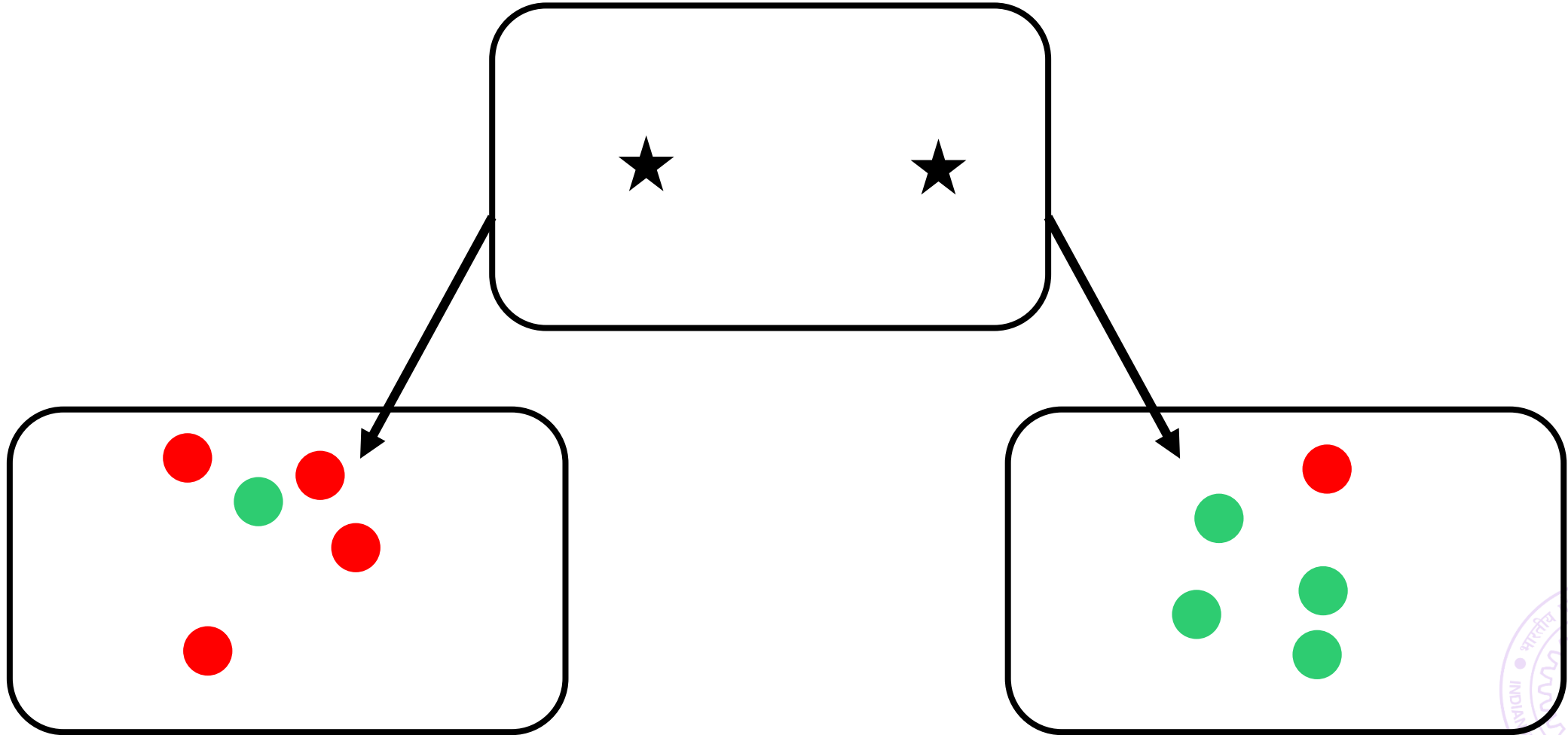
# Node splitting via LwP/linear classifiers 115



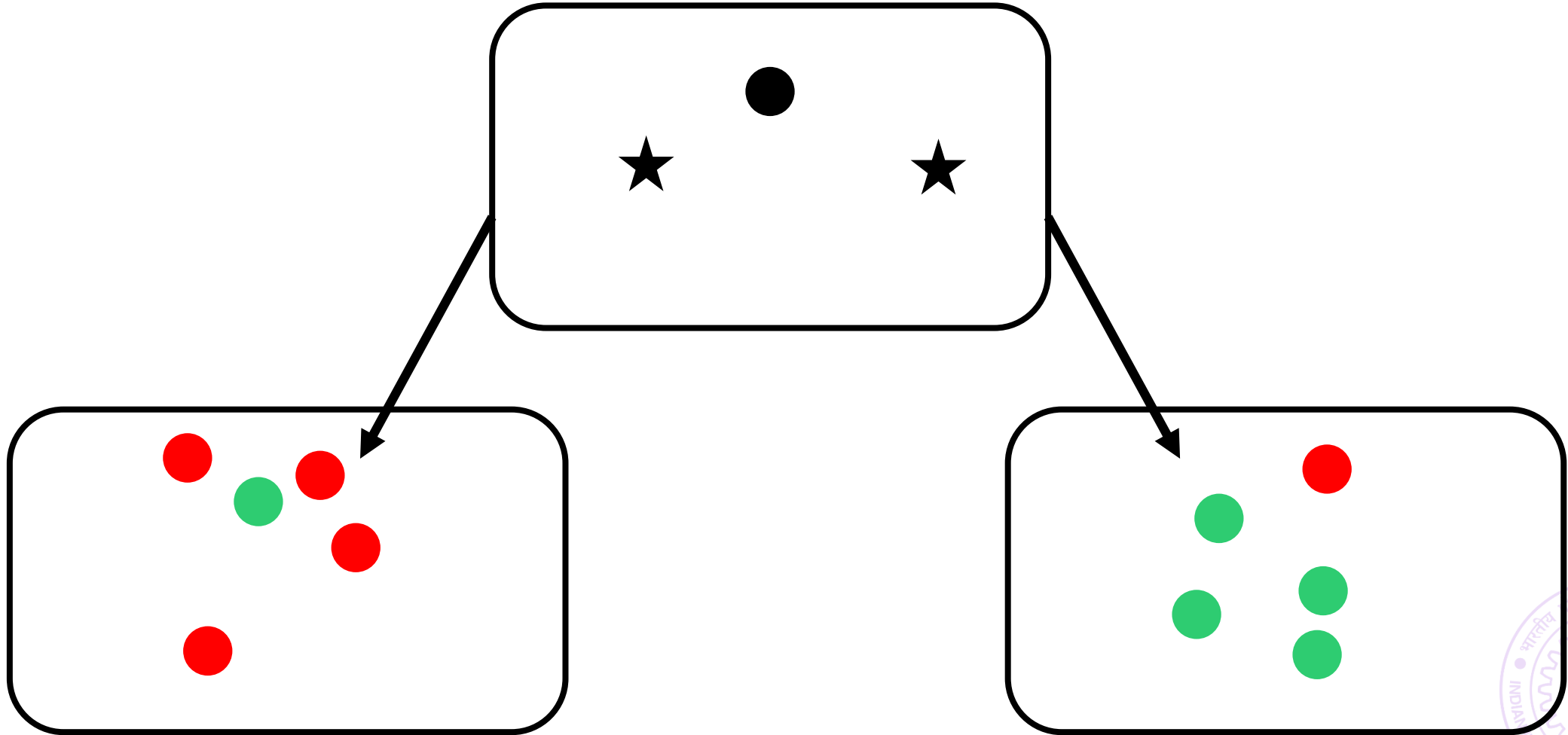
# Node splitting via LwP/linear classifiers 115



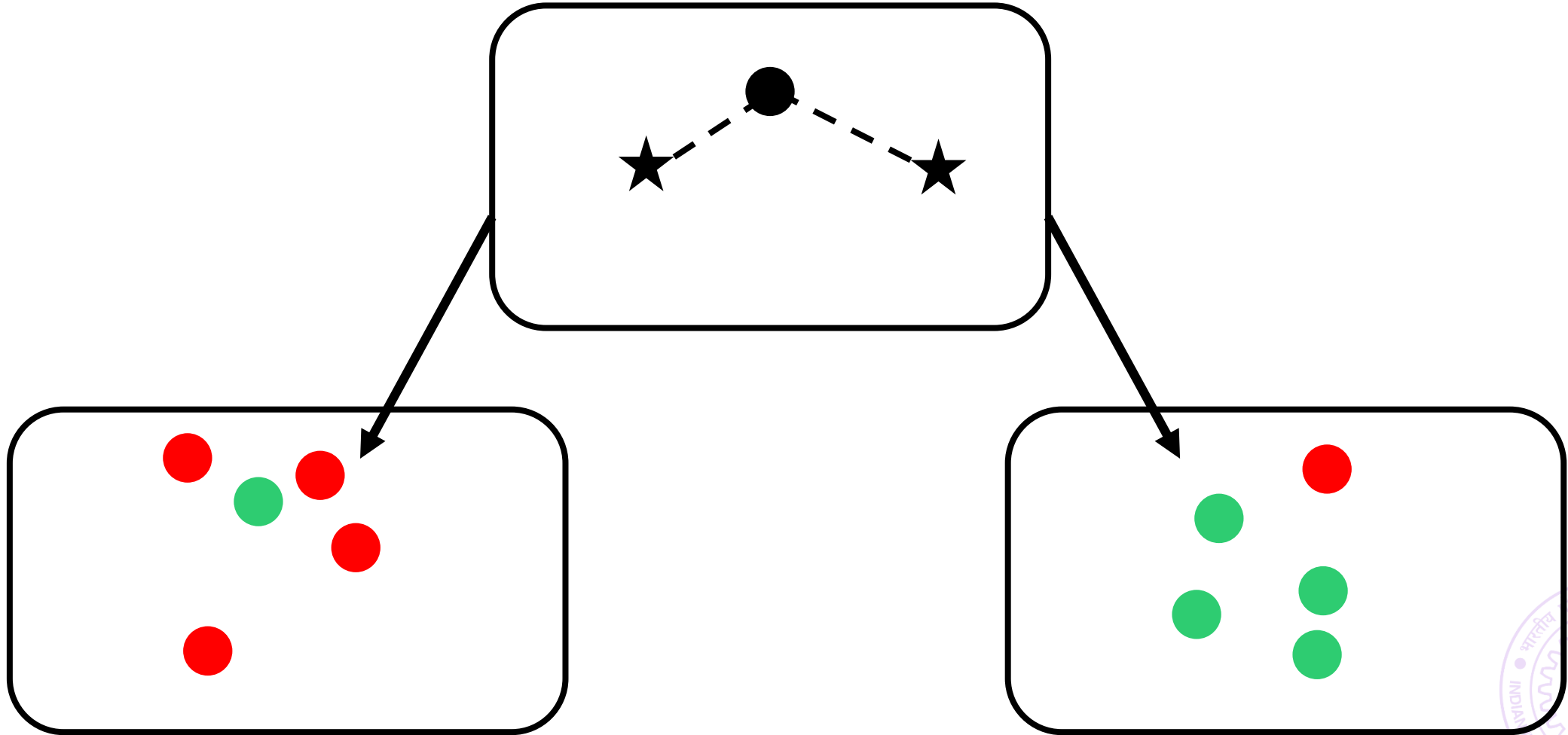
# Node splitting via LwP/linear classifiers 115



# Node splitting via LwP/linear classifiers 115

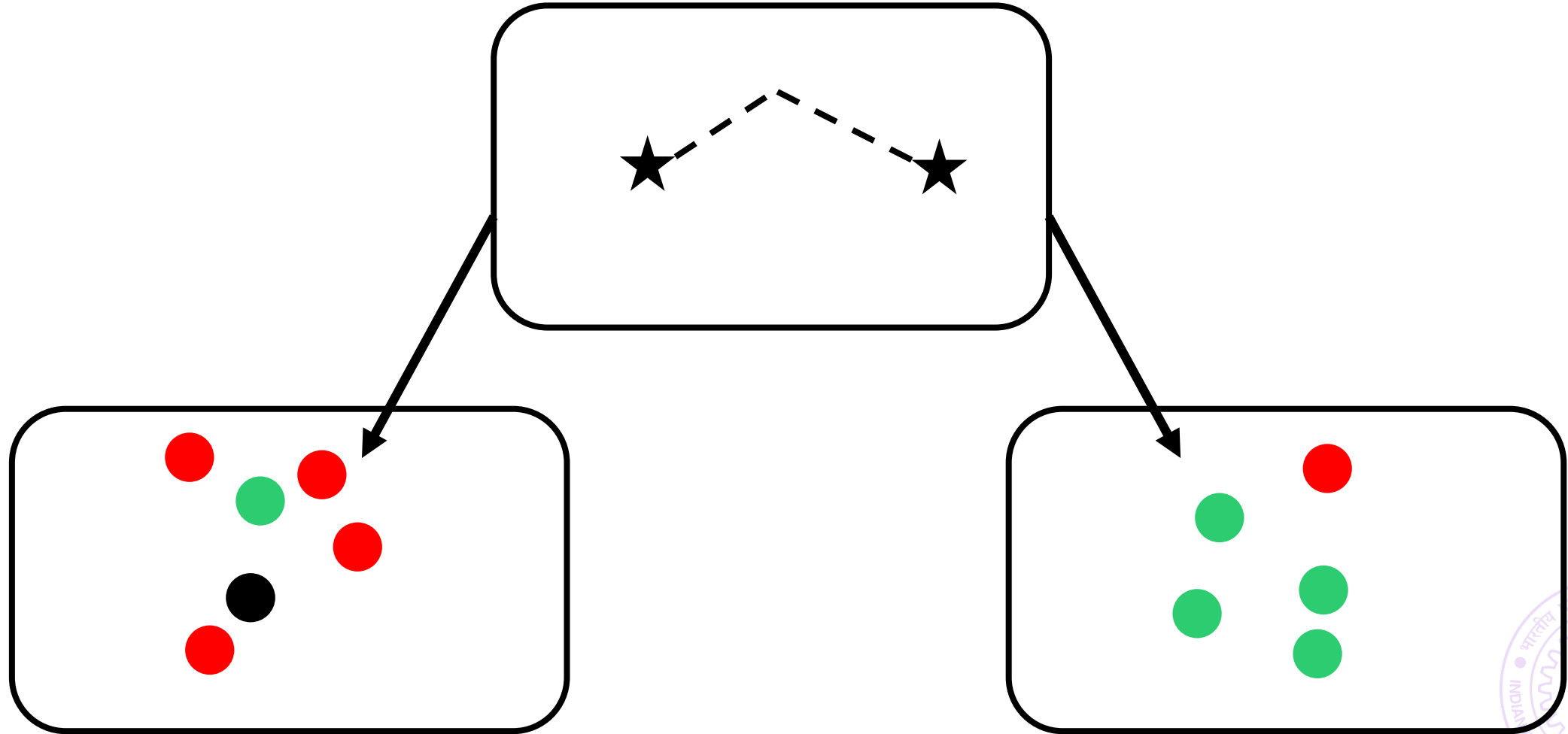


# Node splitting via LwP/linear classifiers 115





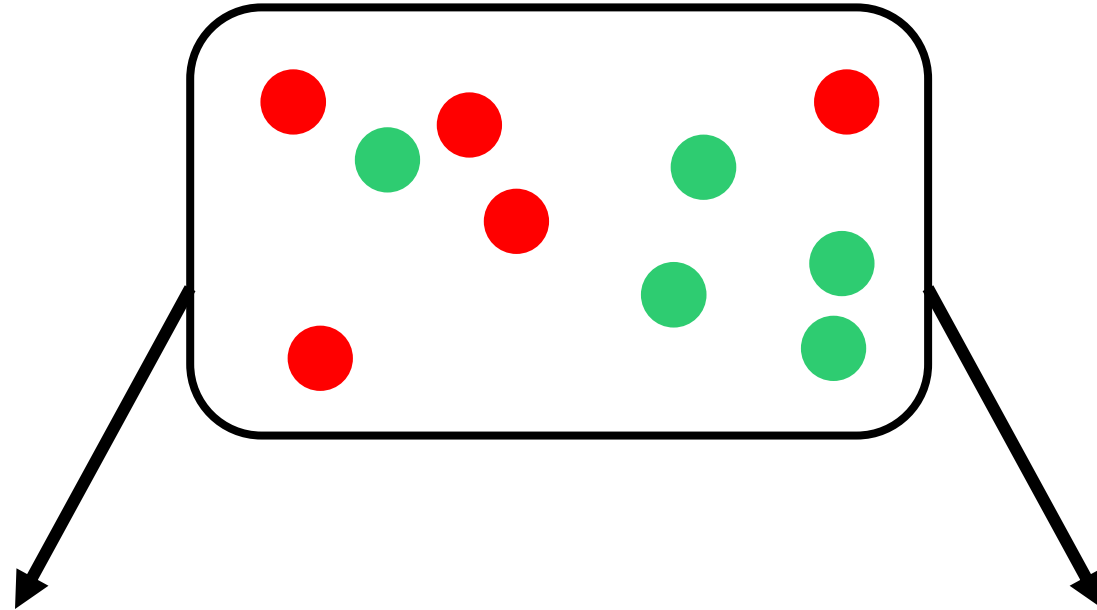
# Node splitting via LwP/linear classifiers 115



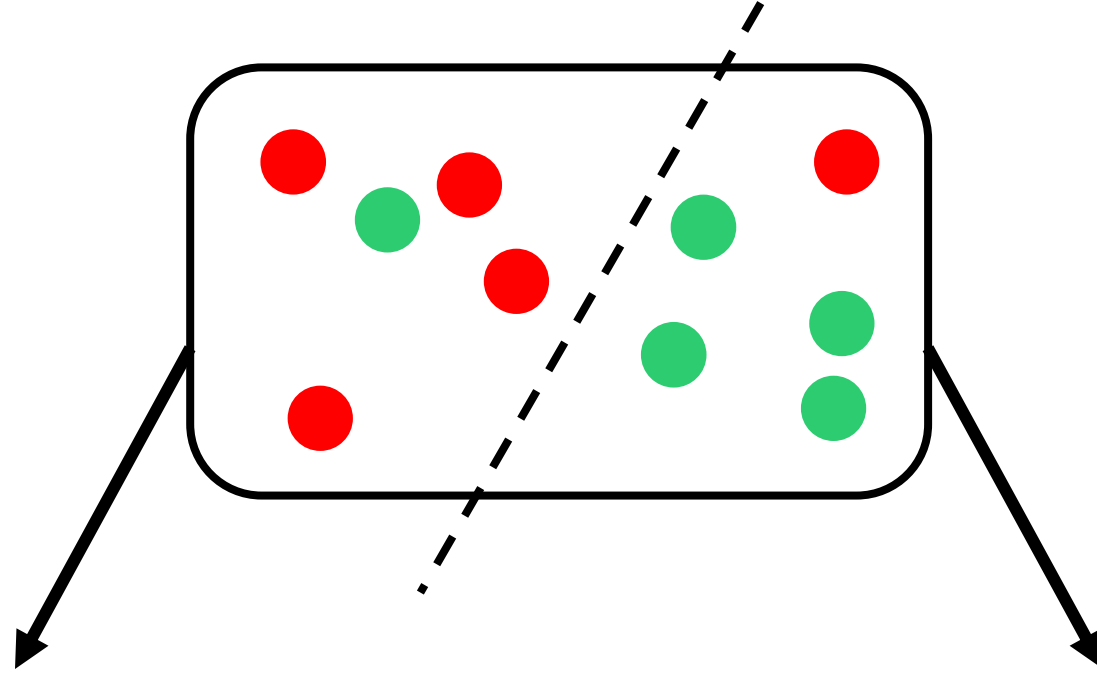
# Node splitting via LwP/linear classifiers 122



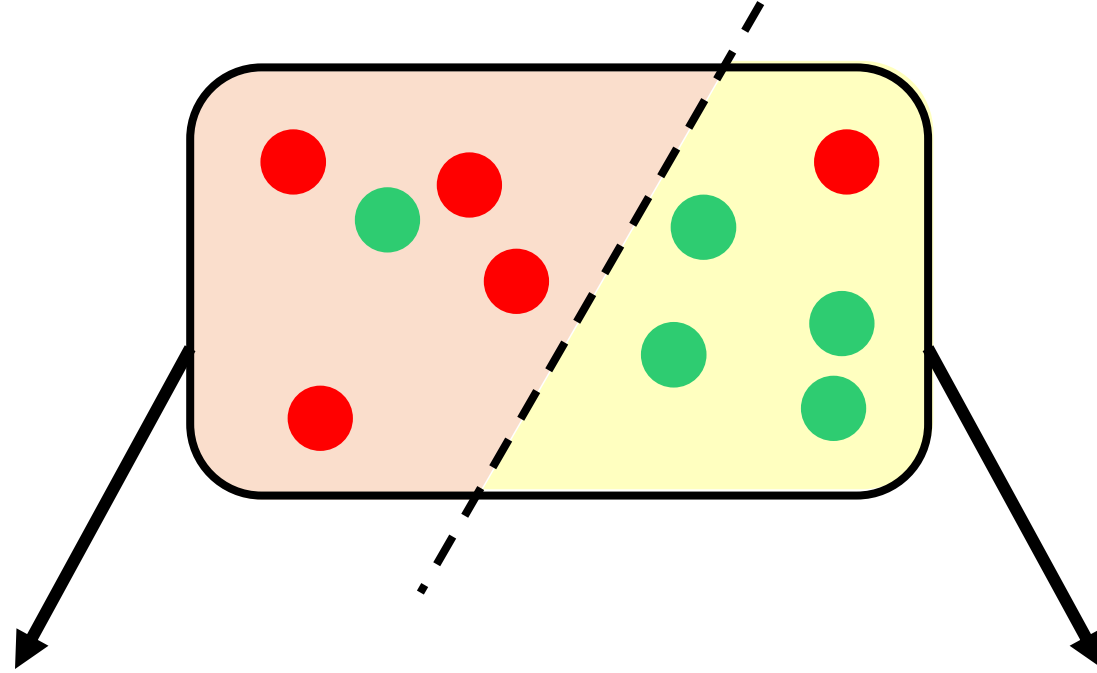
# Node splitting via LwP/linear classifiers 122



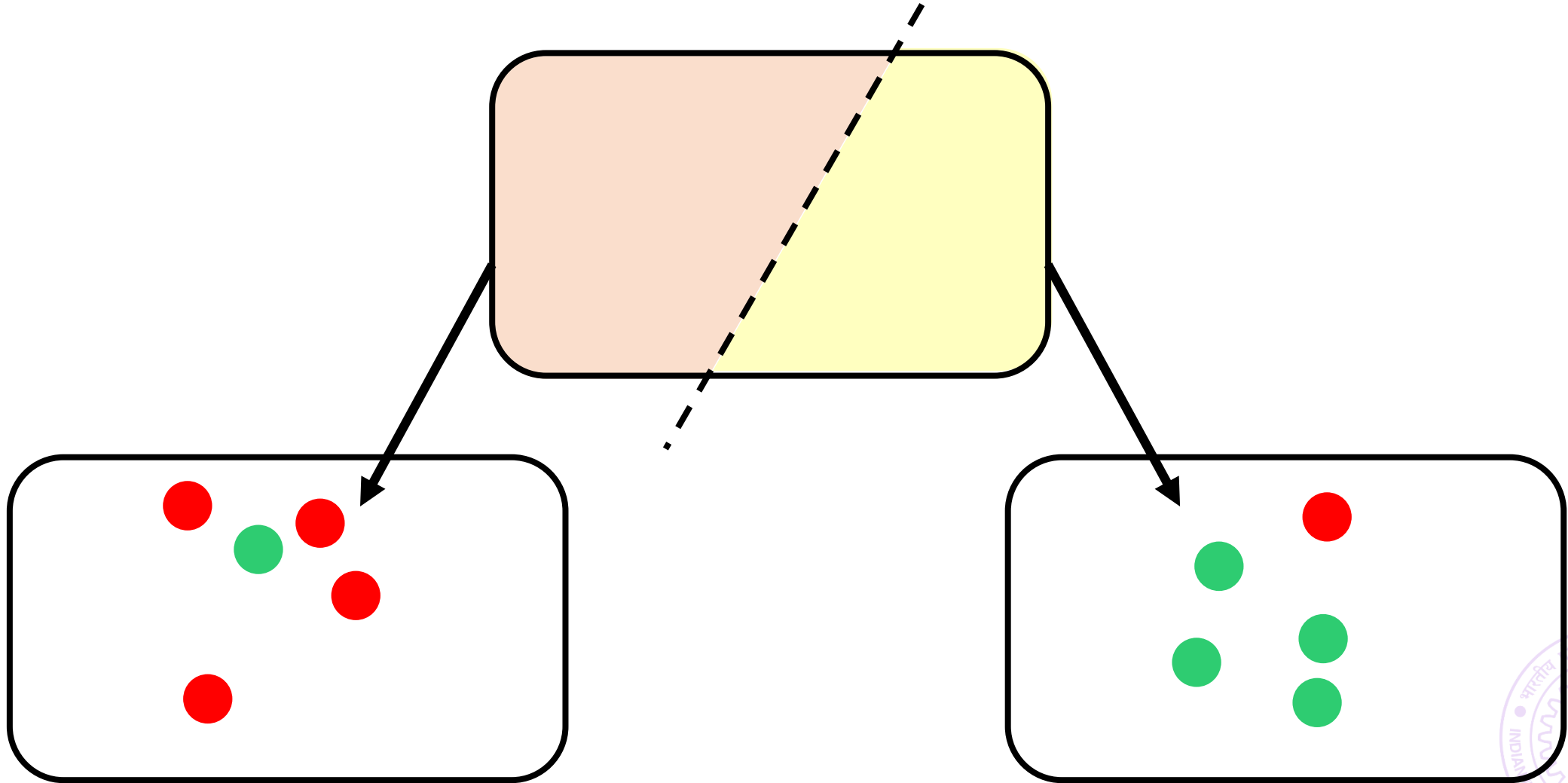
# Node splitting via LwP/linear classifiers 122



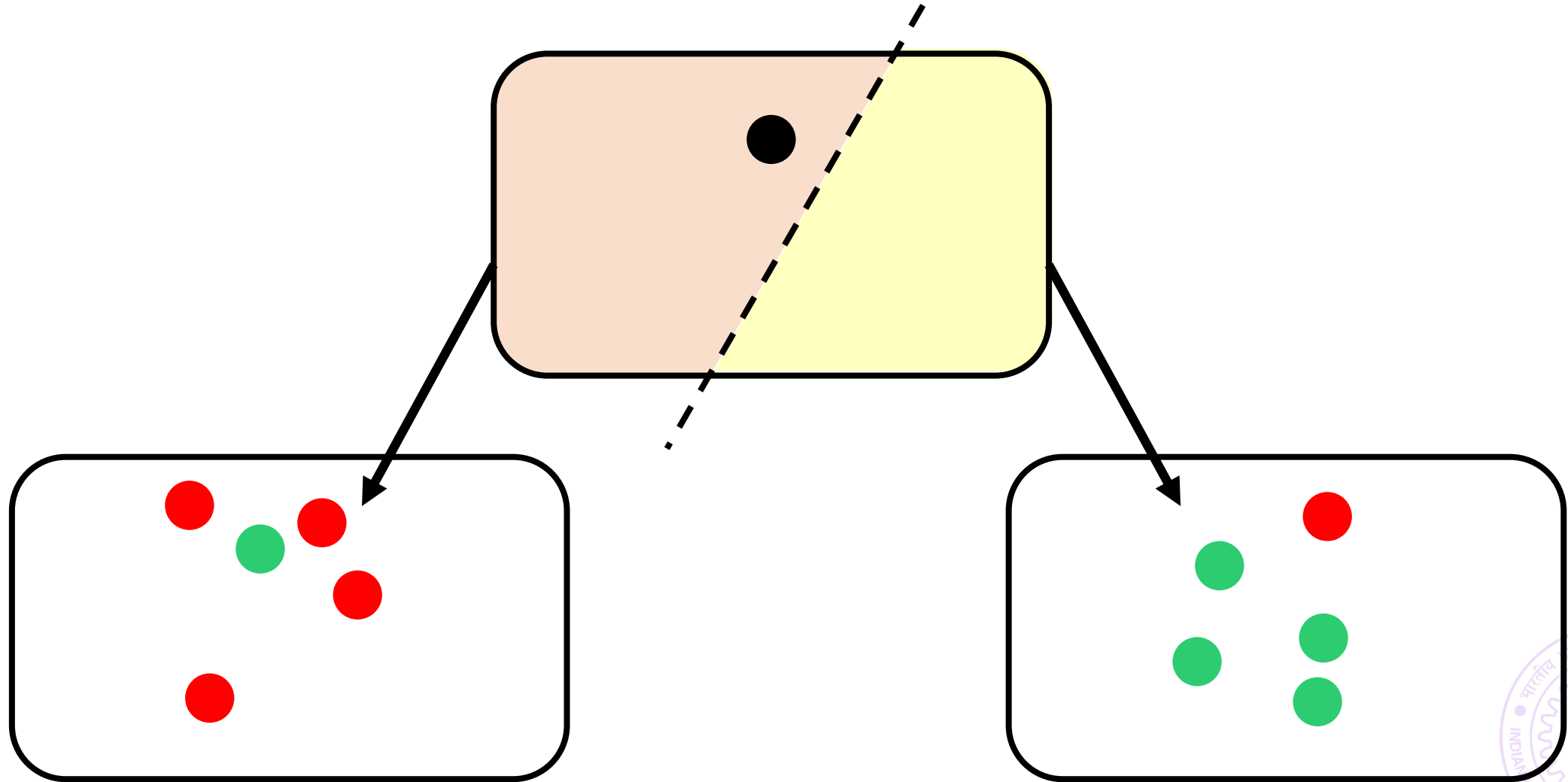
# Node splitting via LwP/linear classifiers 122



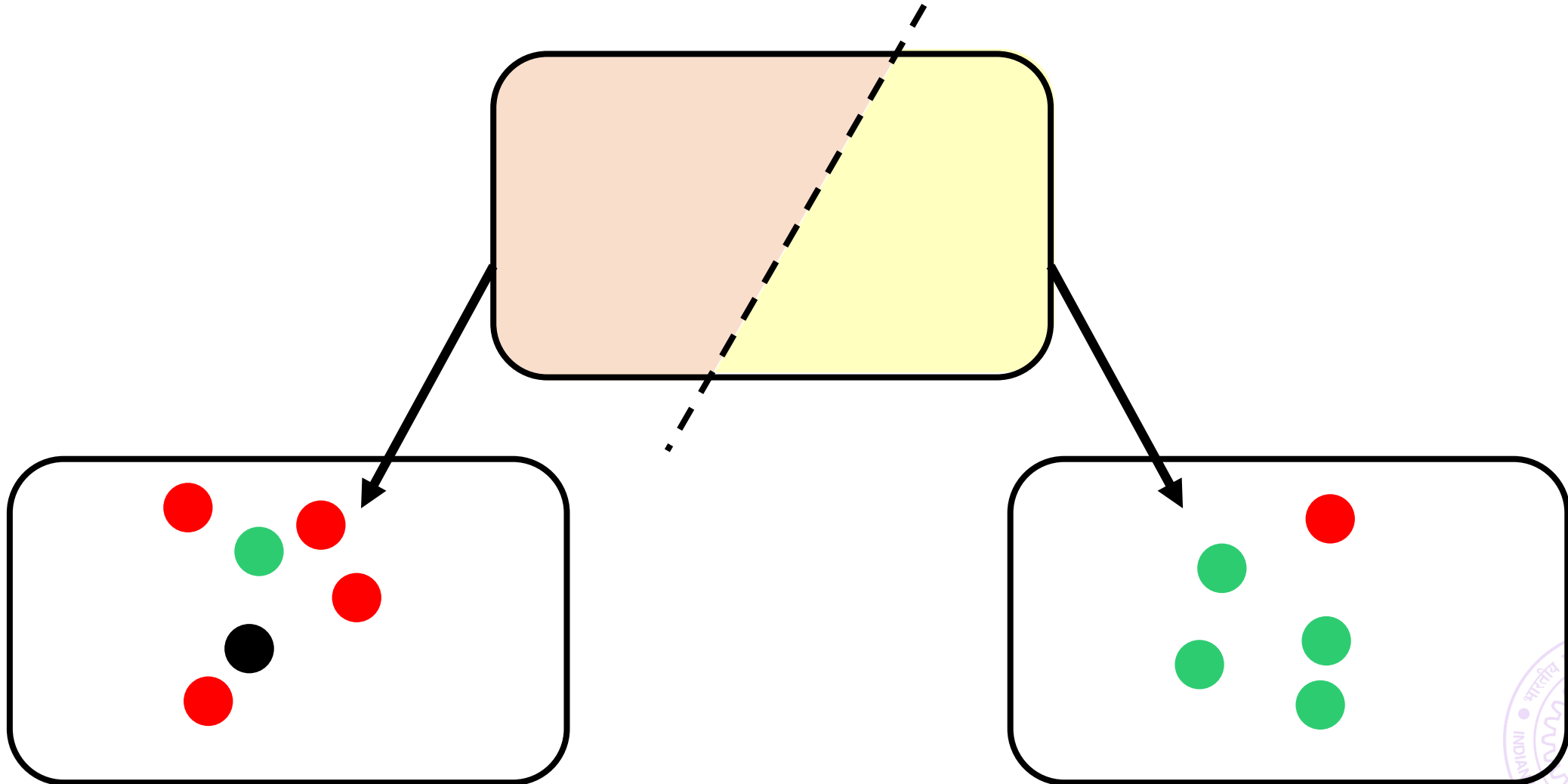
# Node splitting via LwP/linear classifiers 122



# Node splitting via LwP/linear classifiers 122



# Node splitting via LwP/linear classifiers 122





# One Final Recap

129



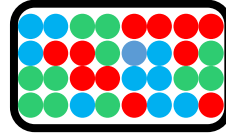
# One Final Recap

129



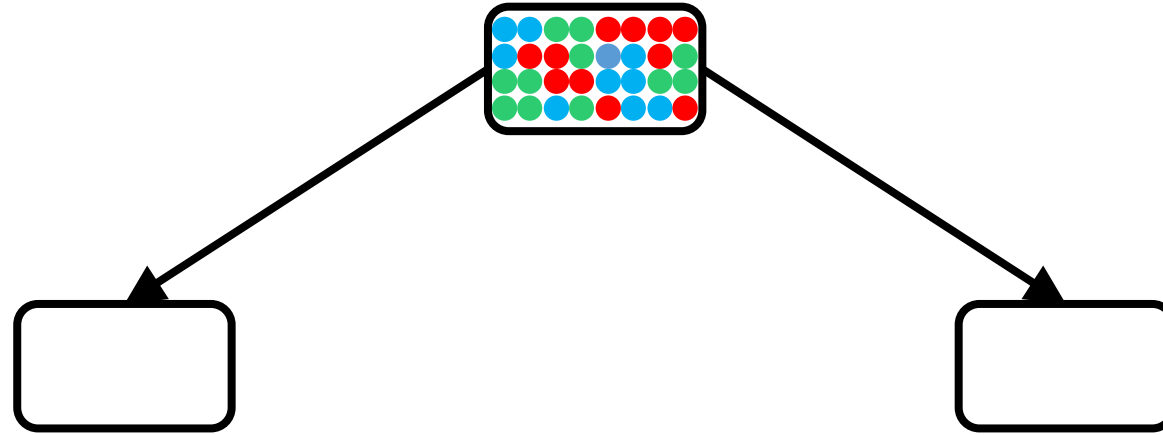
# One Final Recap

129



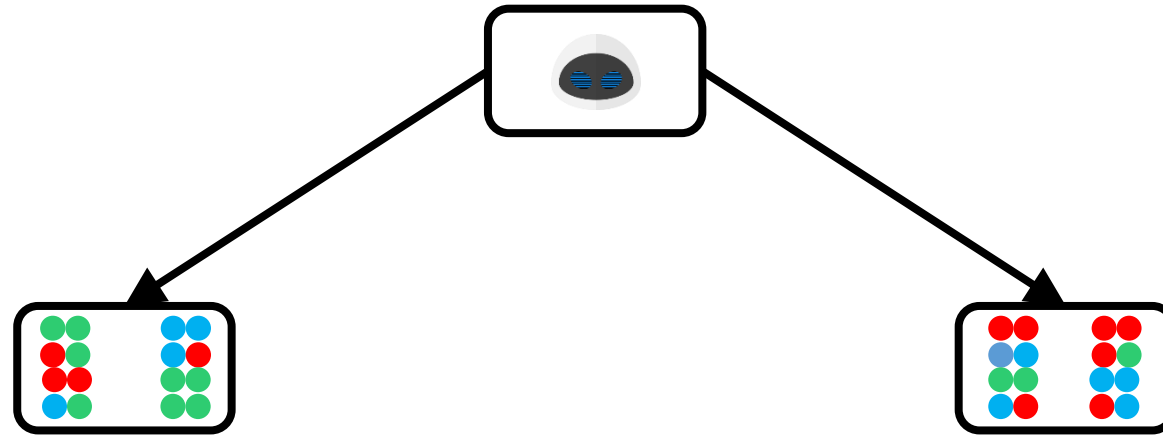
# One Final Recap

129



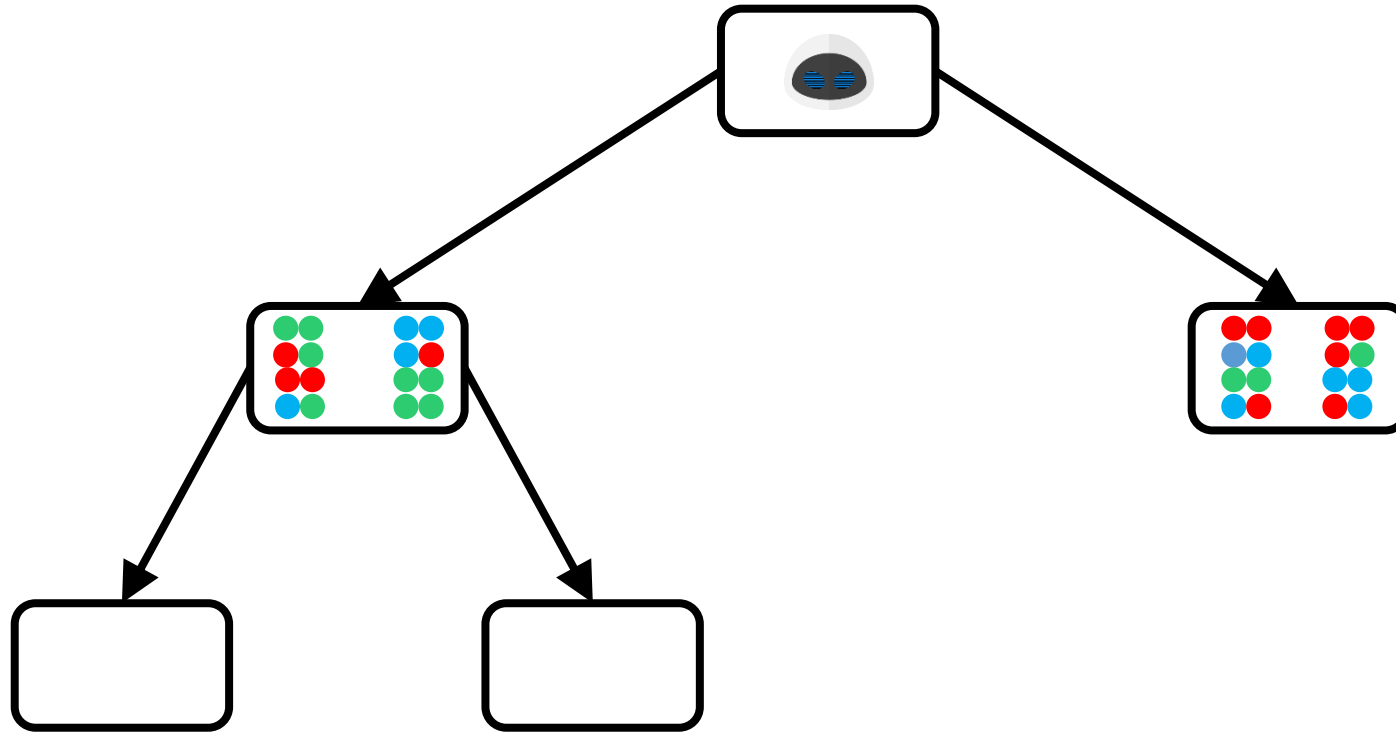
# One Final Recap

129



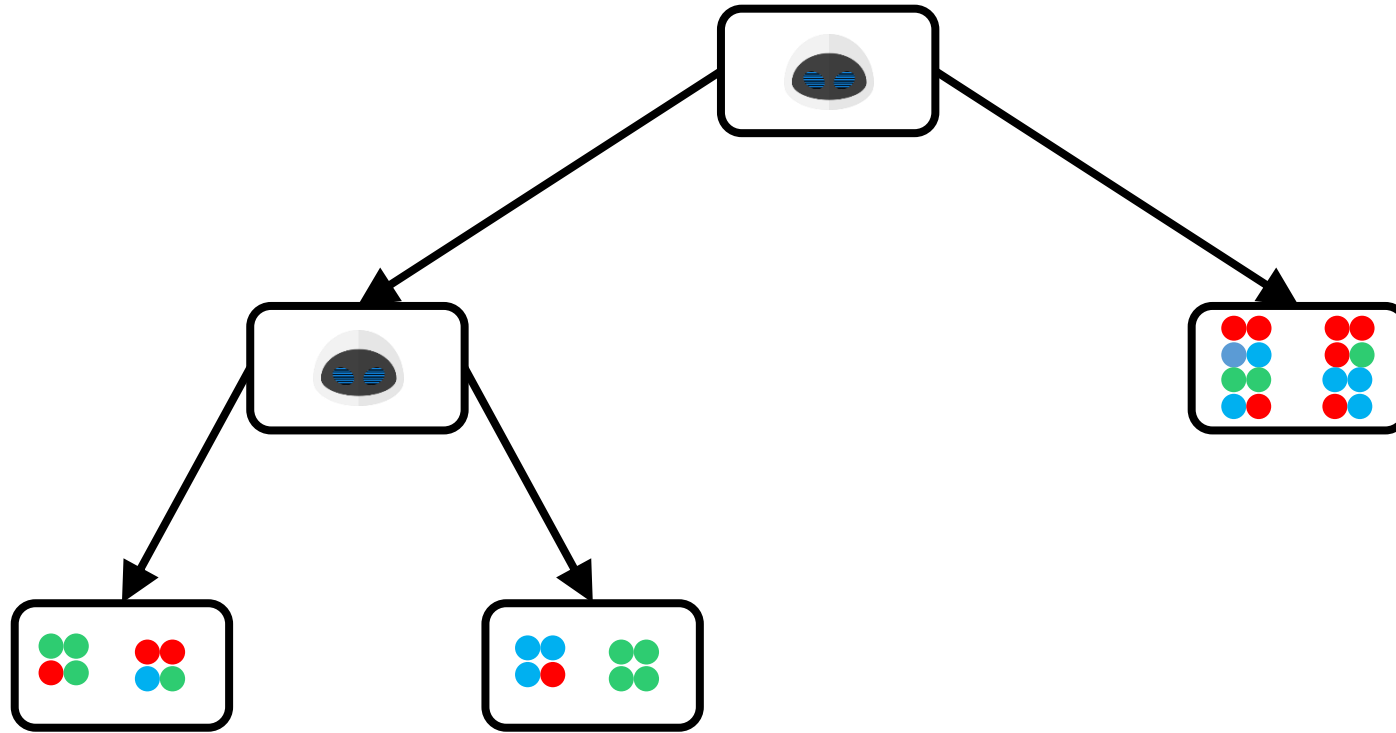
# One Final Recap

129



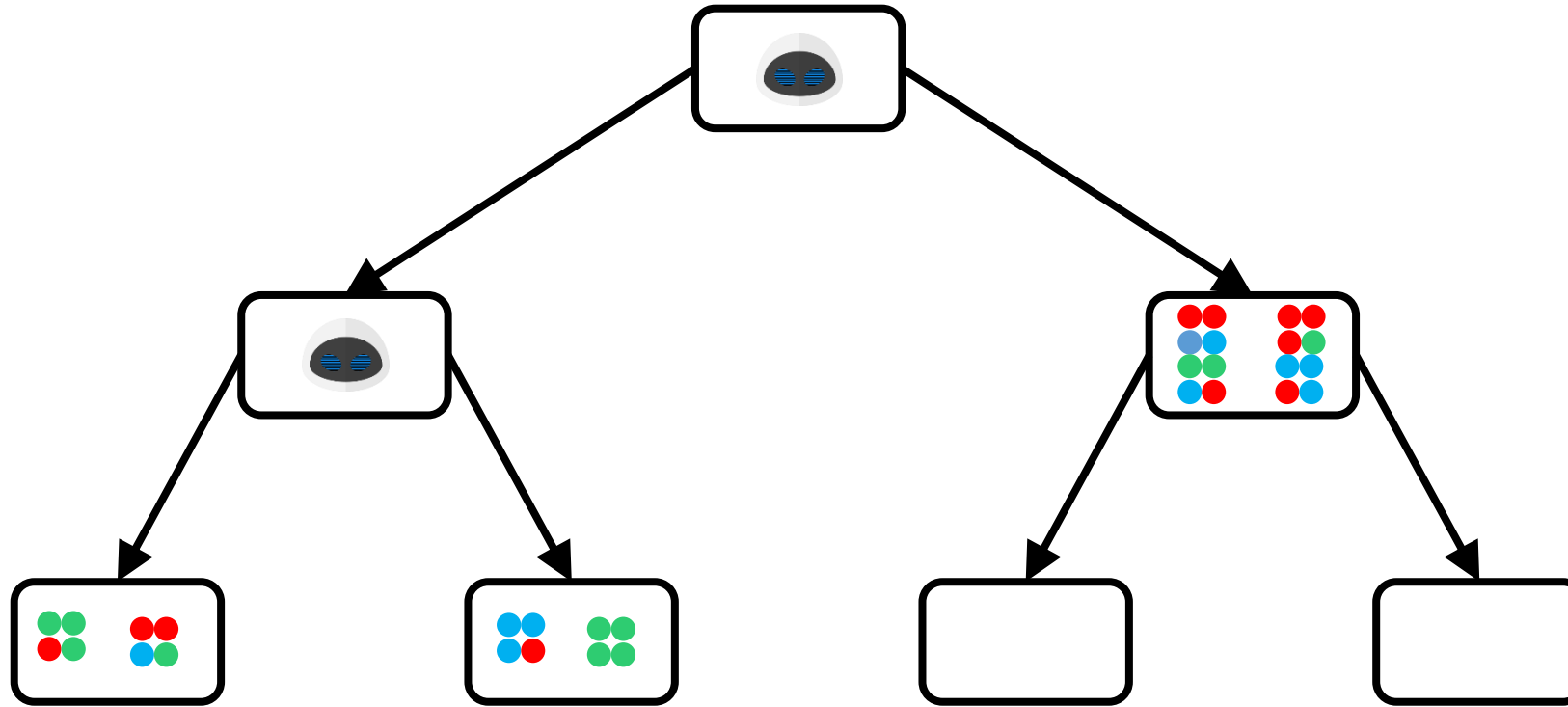
# One Final Recap

129



# One Final Recap

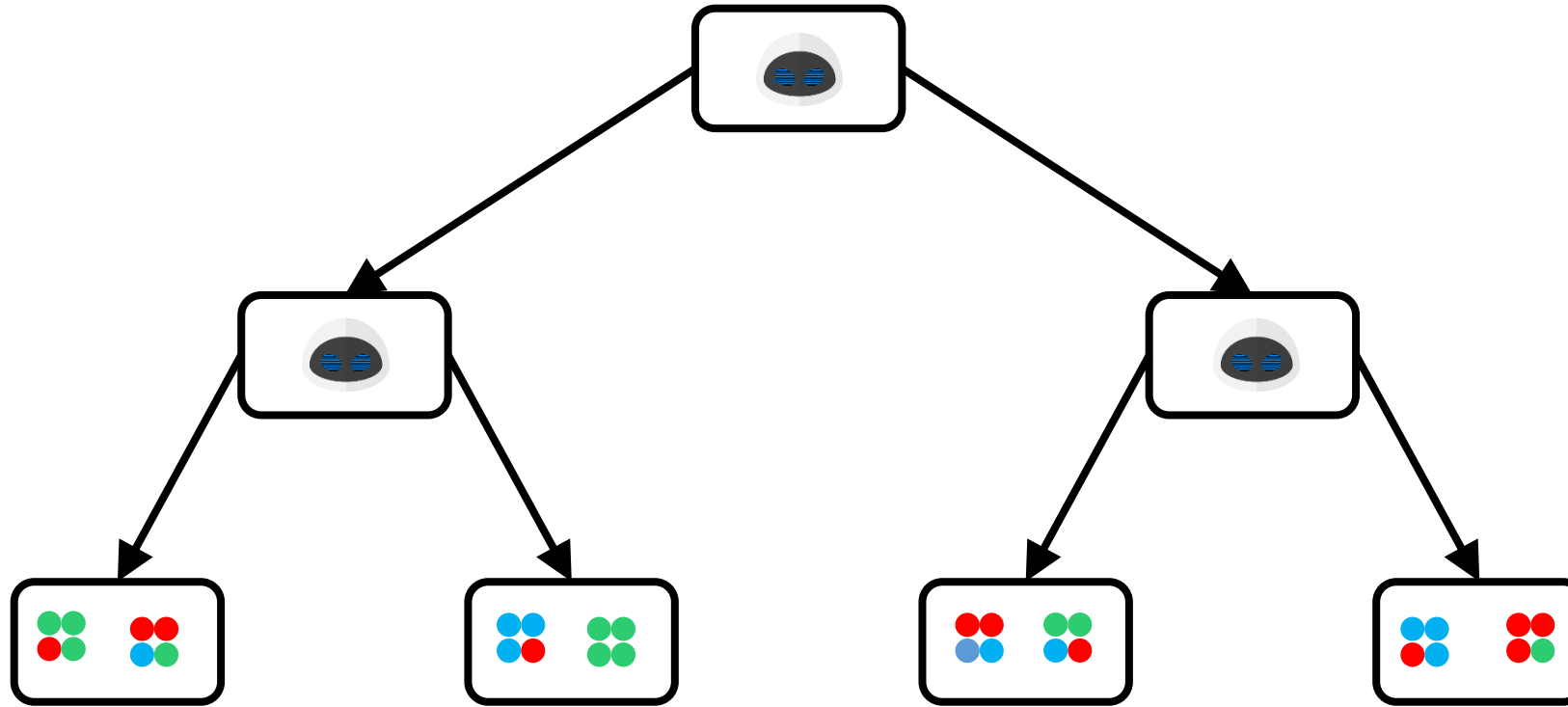
129





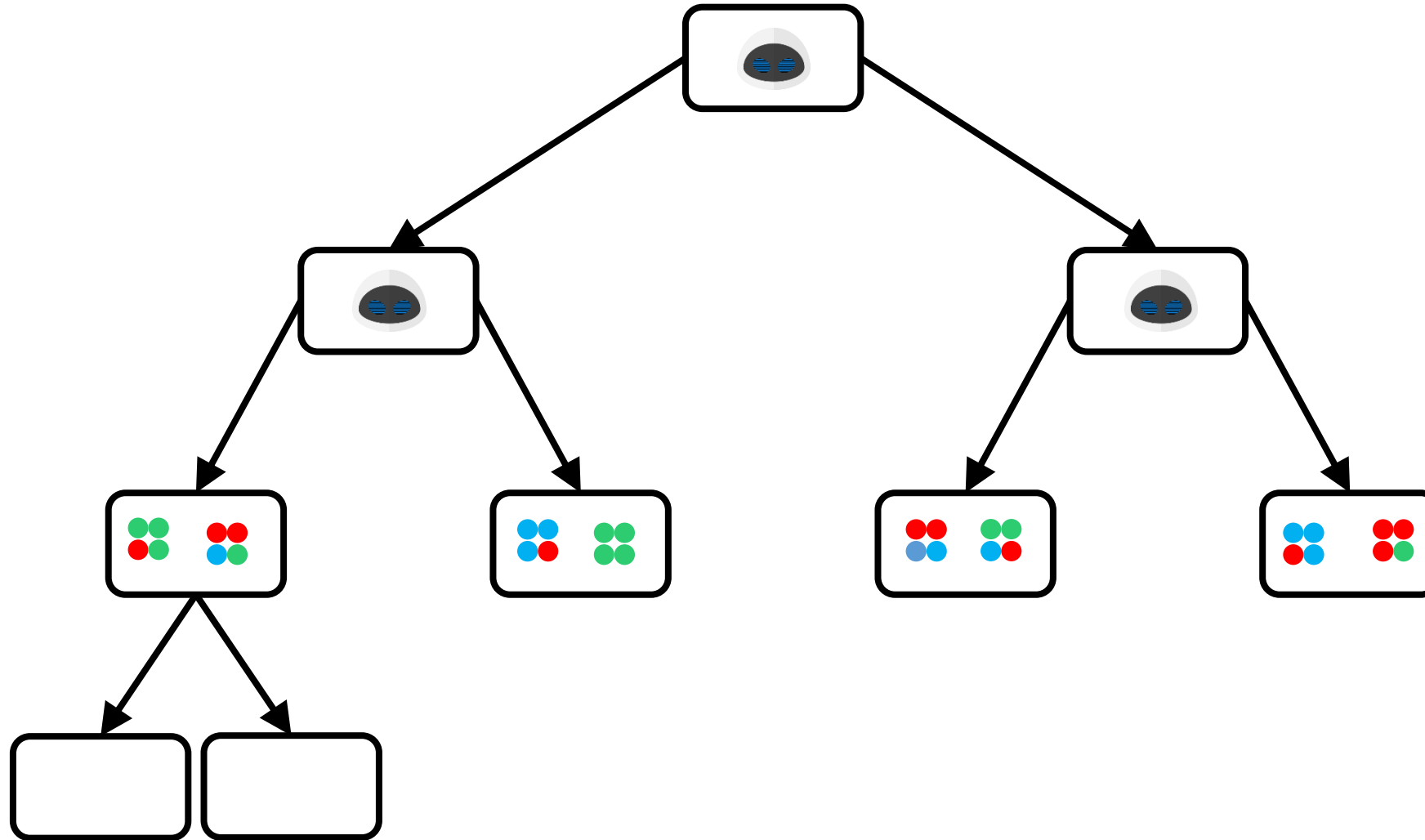
# One Final Recap

129



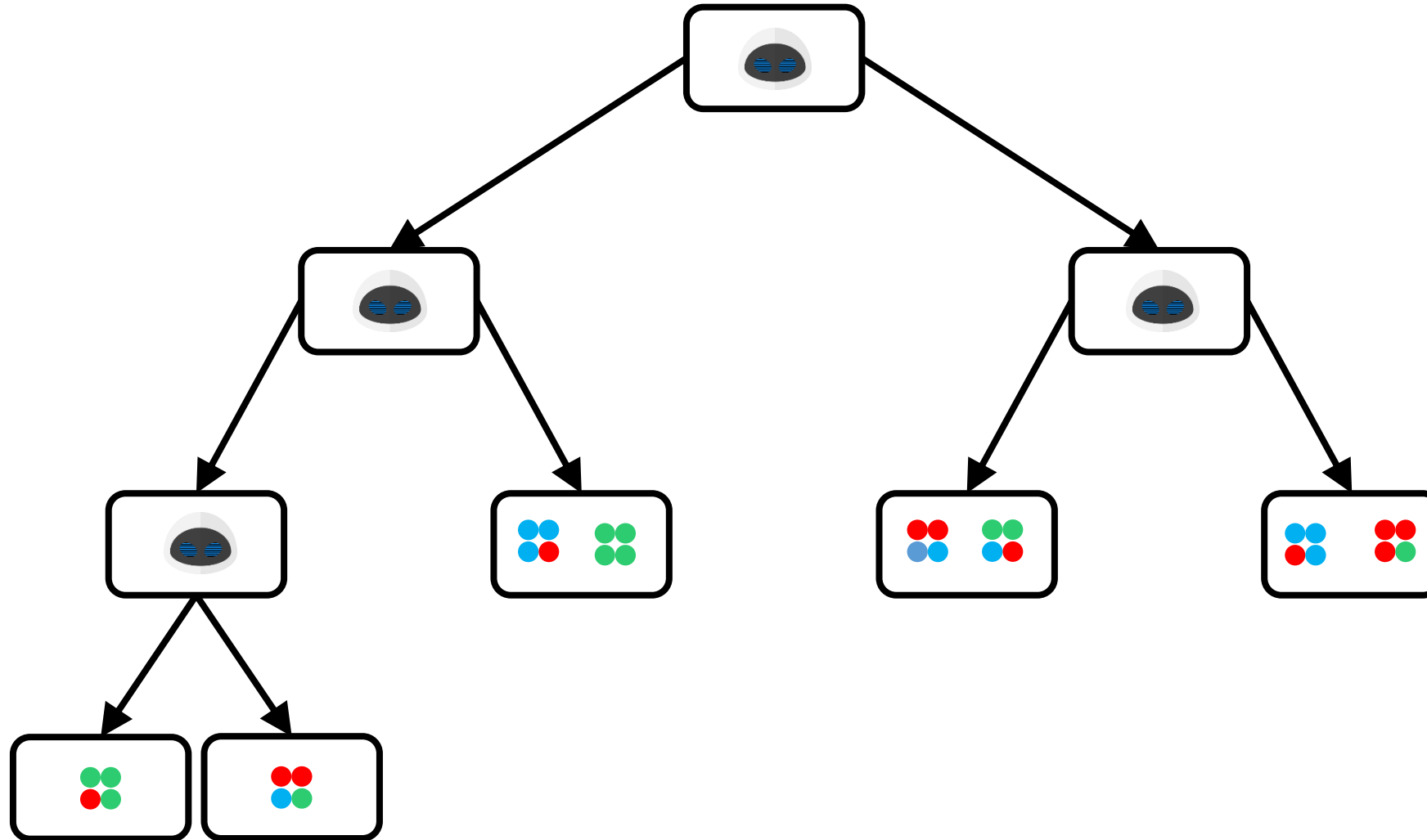
# One Final Recap

129



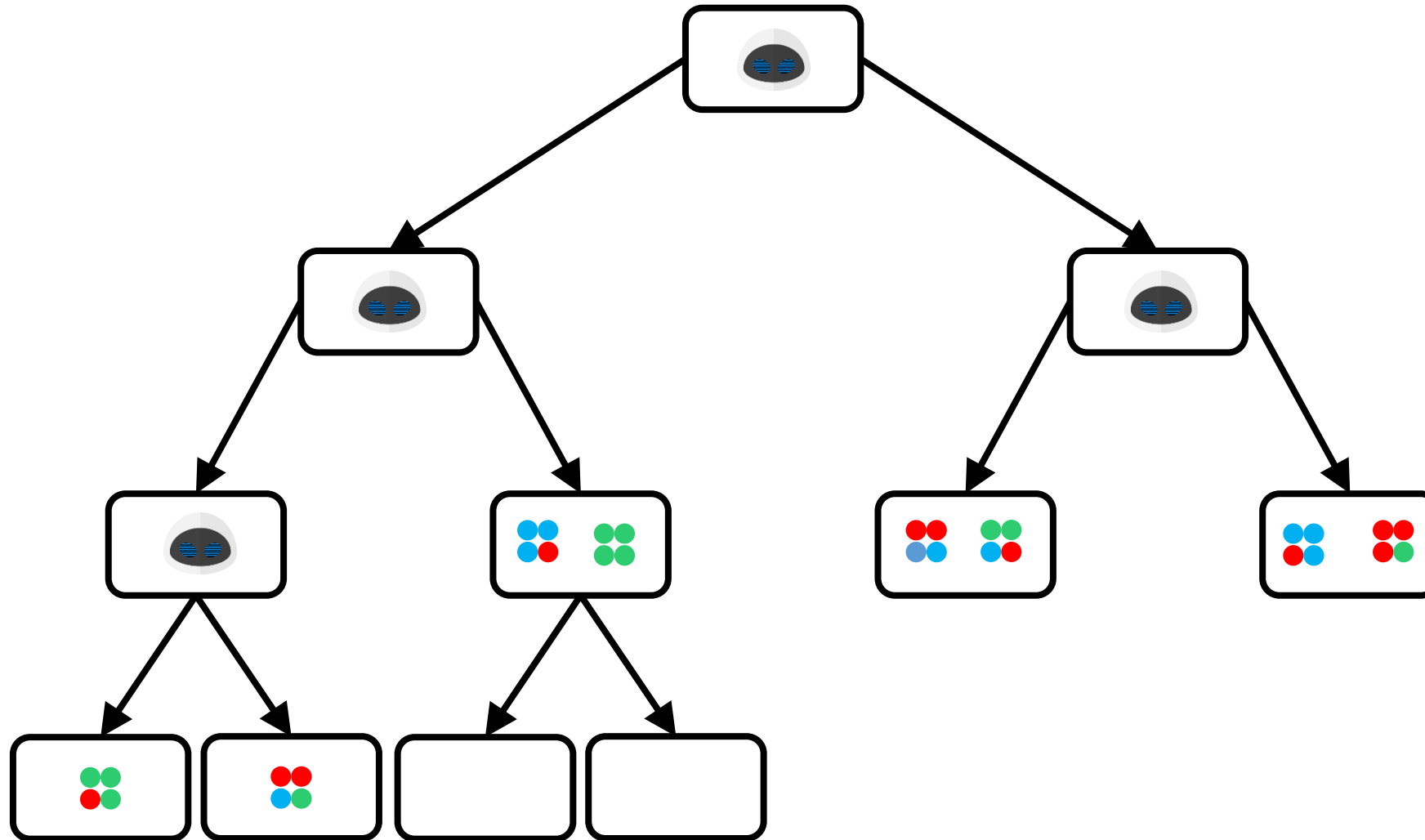
# One Final Recap

129



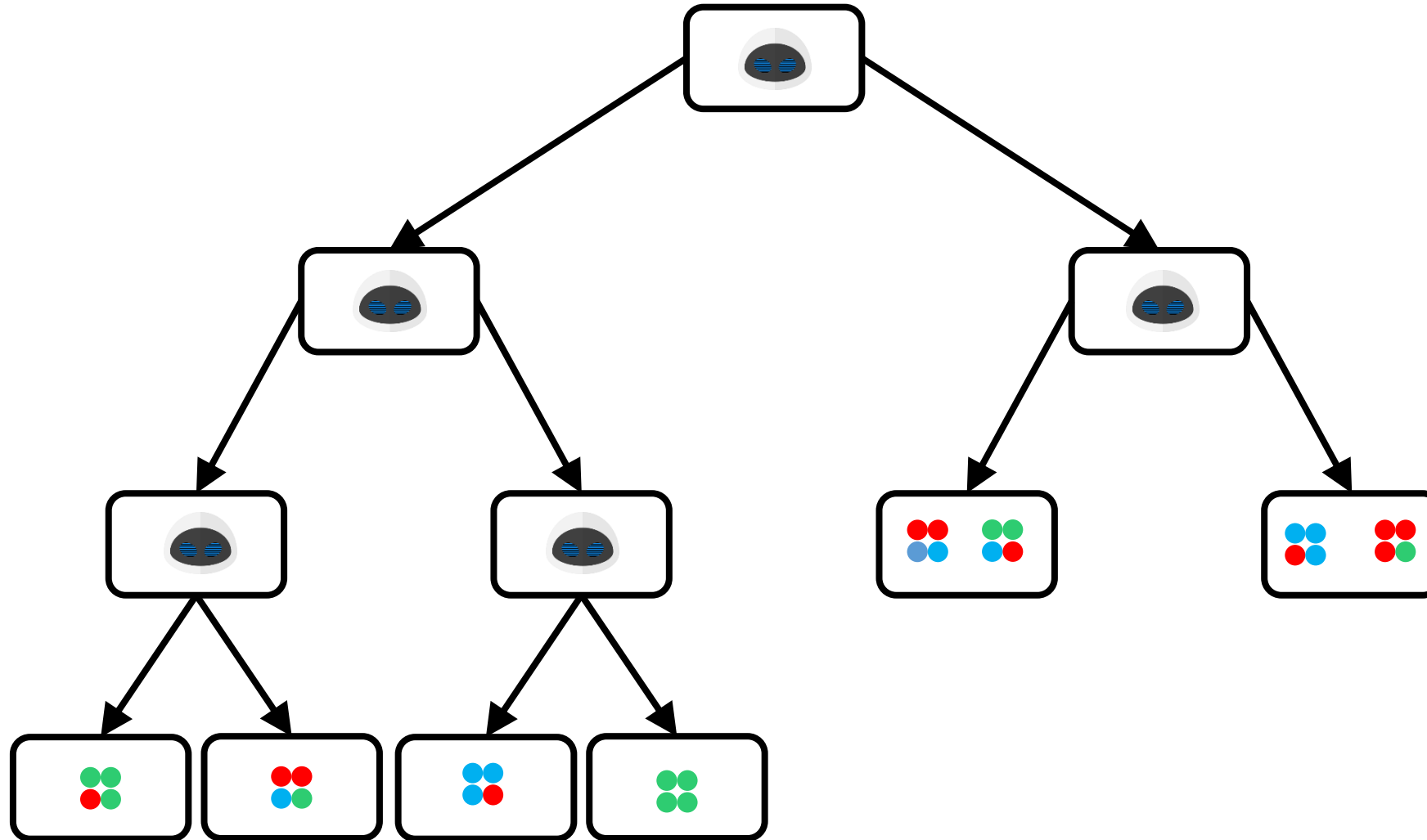
# One Final Recap

129



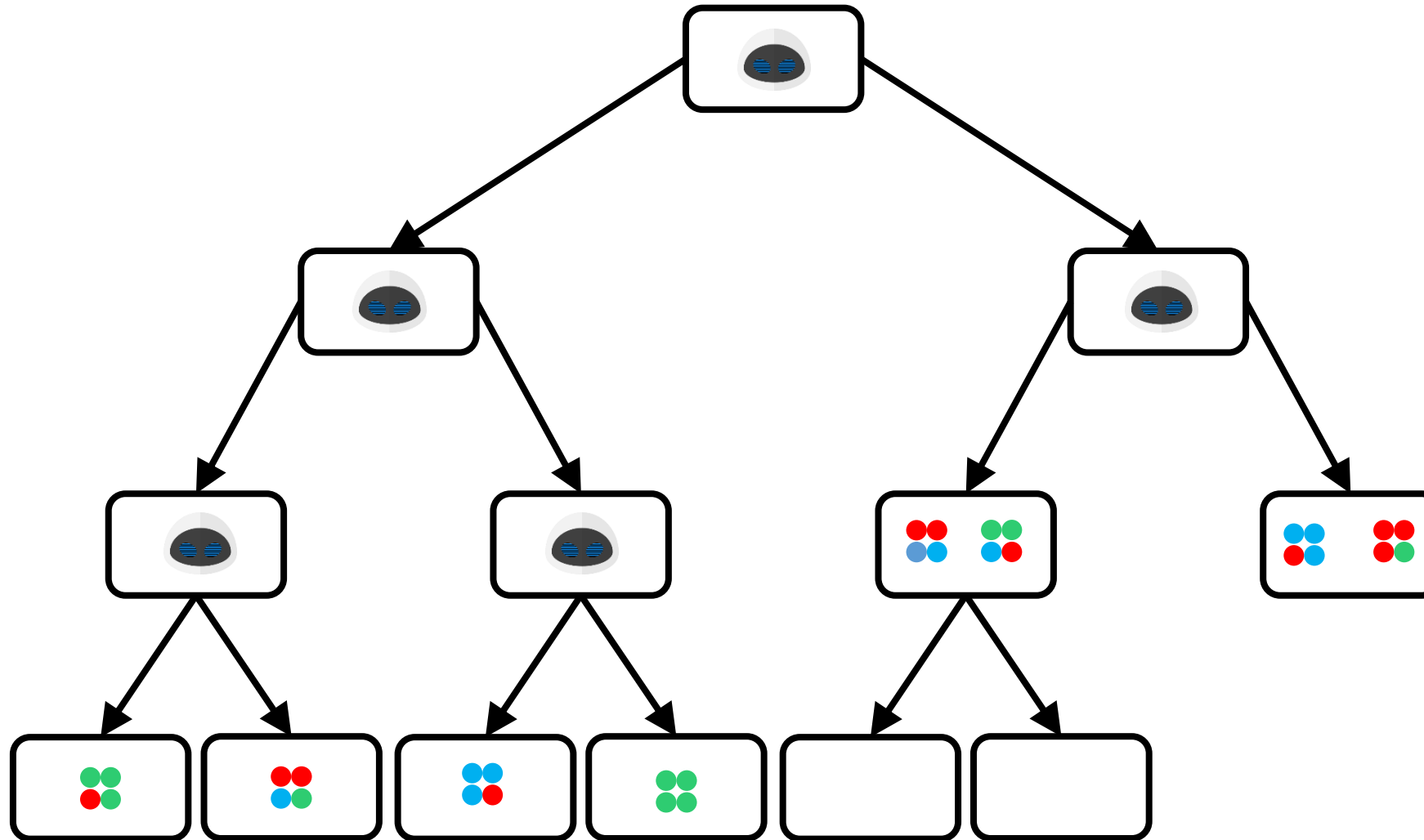
# One Final Recap

129



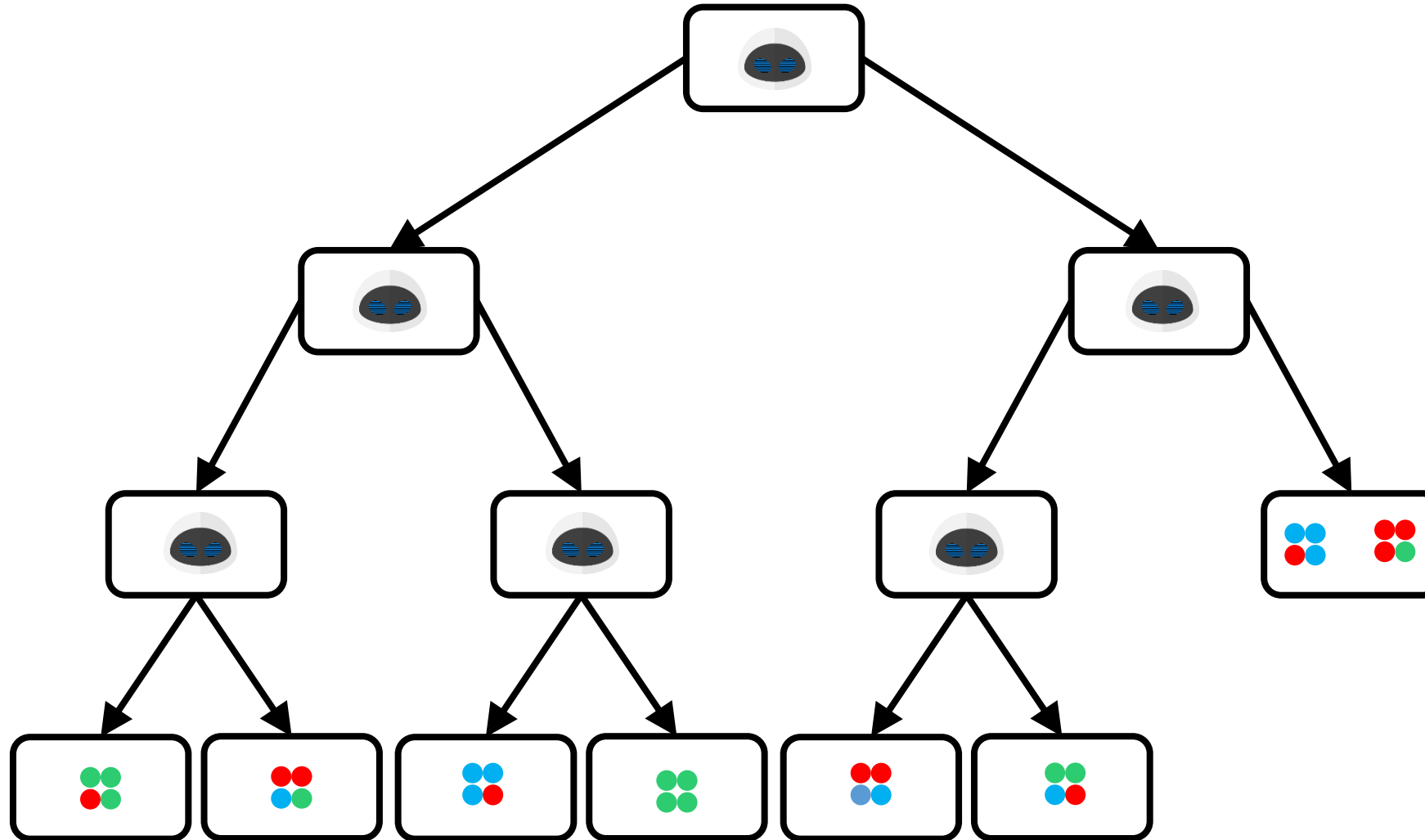
# One Final Recap

129



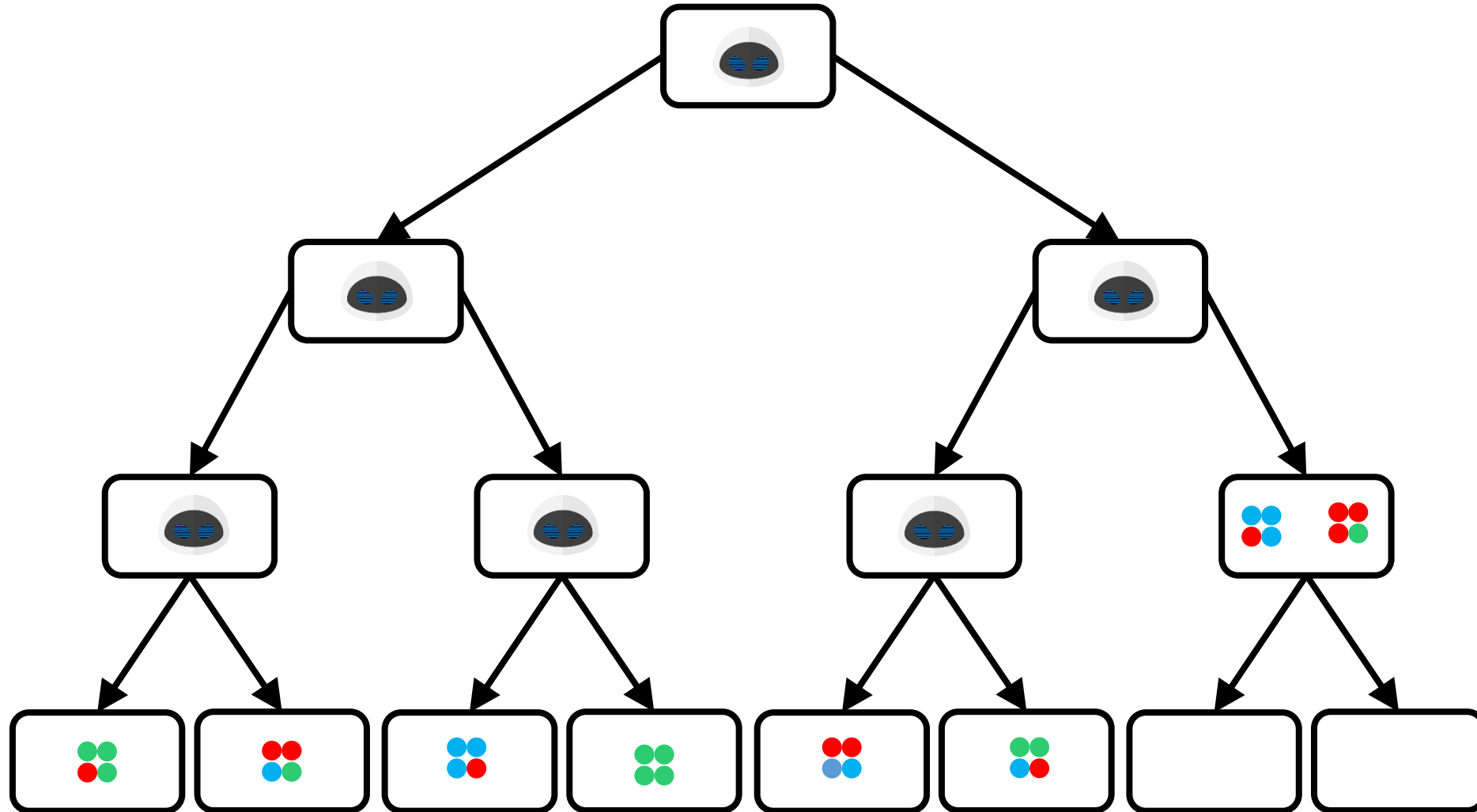
# One Final Recap

129



# One Final Recap

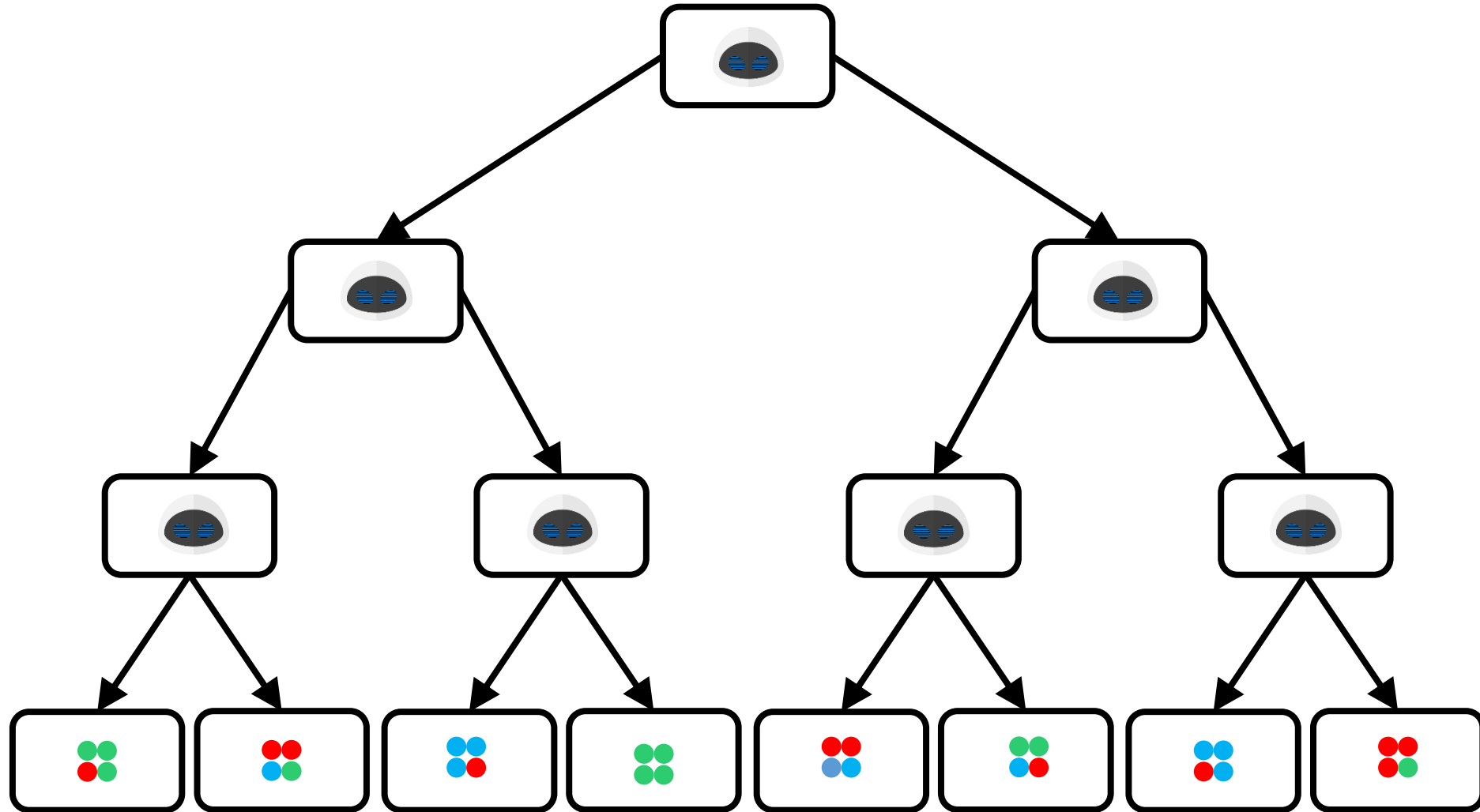
129





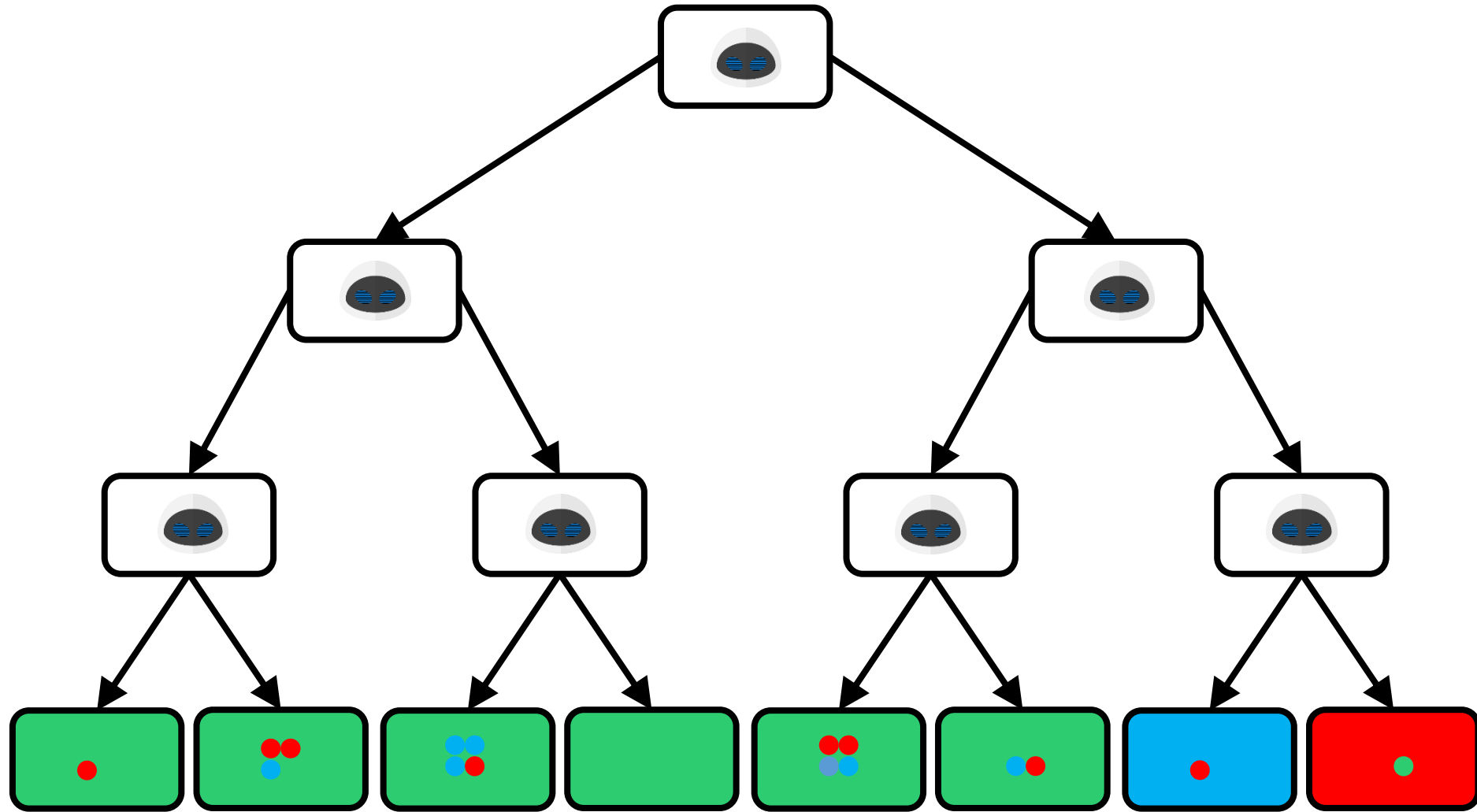
# One Final Recap

129



# One Final Recap

129



# Pruning Strategies

147

Stop if node is pure or almost pure

Stop if all features exhausted – avoid using a feature twice on a path

Limits depth of tree to  $d$  (num of dimensions)

Can stop if a node is ill-populated i.e. has few training points

Can also (over) grow a tree and then merge nodes to shrink it

Merge two leaves and see if it worsens performance on the validation set or not – rinse and repeat

Use a validation set to make these decisions (never touch test set)



# Decision Trees - Lessons

148

Very fast at making predictions (if tree is reasonably balanced)

Can handle discrete data (even non numeric data) as well – e.g. can have a stump as: blood group AB or O go left, else go right

LwP, NN have difficulty with such non-numeric and discrete data since difficult to define distance and averages with them (however, there are workarounds to do NN/LwP with discrete data as well)

Tons of DT algorithms – both classical (ID3, C4.5) as well as very recent (GBDT, LPSR, Parabel) – DTs are versatile and very useful

Reason: DT learning is an NP hard problem – no single algo ☹️

If you think you have a better way of splitting nodes or handling leaf nodes, it might be the next big thing in DT learning 😊

