

# Decision Trees

# Decision Trees - 1

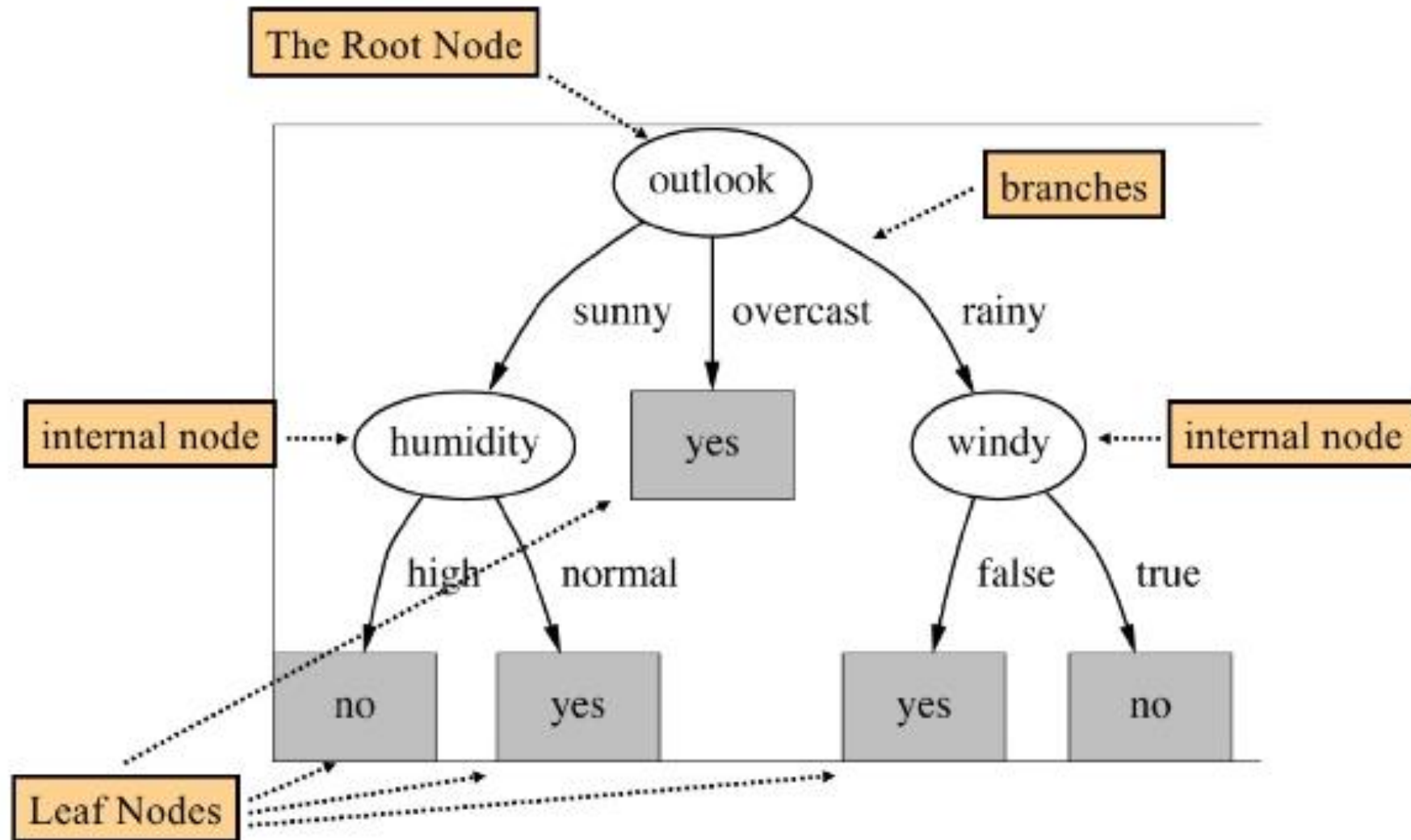
- Algorithms used for both Classification and Regression
- Works effectively with non-linear dataset
- DT can be looked at like “rules” that can be understood by humans and implemented on datasets
- Uses the Greedy algorithm technique
- Core algorithm to build decisions is called **ID3**, that employs a top-down approach
- Dataset is split on the ***most significant feature*** (using ***Entropy / Information Gain***)
- DT represents an inverted tree having the following attributes:
  - ❑ **Decision node** : Test for split of an attribute
  - ❑ **Edge** : split of an attribute
  - ❑ **Leaf node** : value of the target attribute
  - ❑ **Path** : a series of test to arrive at the final decision
- Using recursion, sub-trees are formed based on features not used in the higher nodes
- Divide and rule
- DT splits data until it reaches a “pure” state
  - Pure subset is one where there are only **positive** outcomes. No further split

## Greedy algorithm technique

- A choice made which seems appropriate at that point of time
- A local-optimum choice that would lead to a global-optimum solution
  - But doesn't happen always
- Algorithm does not go back and reverse its decision
  - has only 1 shot to make the local optimum choice

**IDE, C4.5, C5.0 -> Entropy and Information Gain**  
**CART -> Gini Index model**

# Decision Tree Terminology

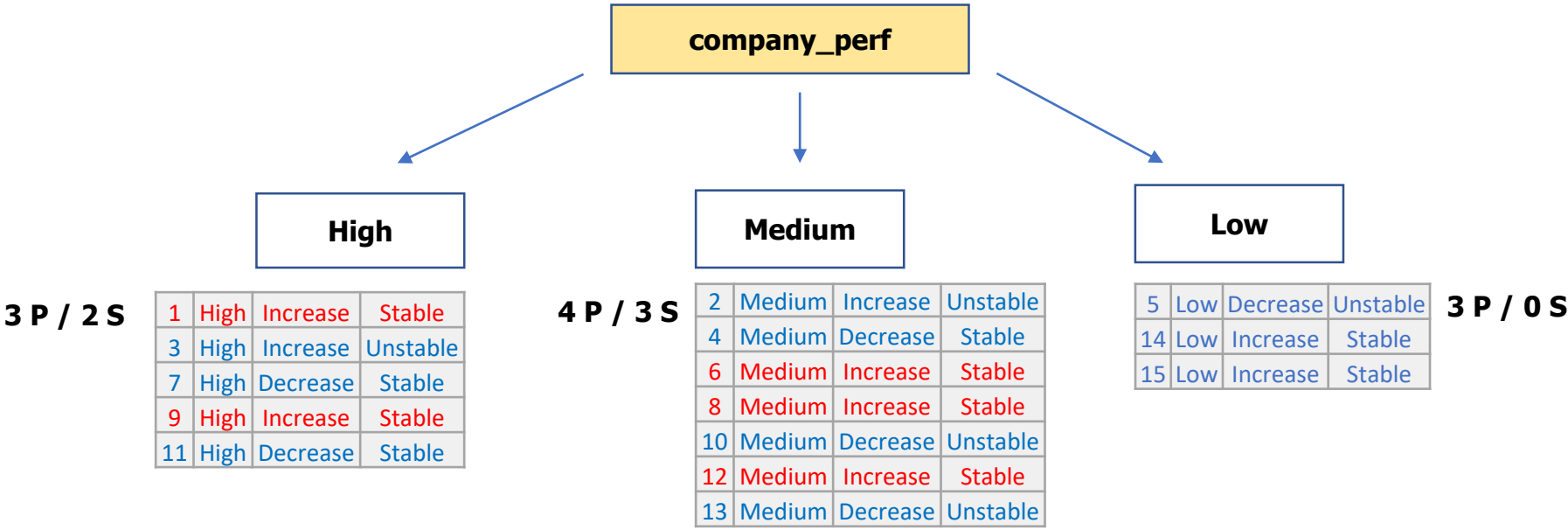


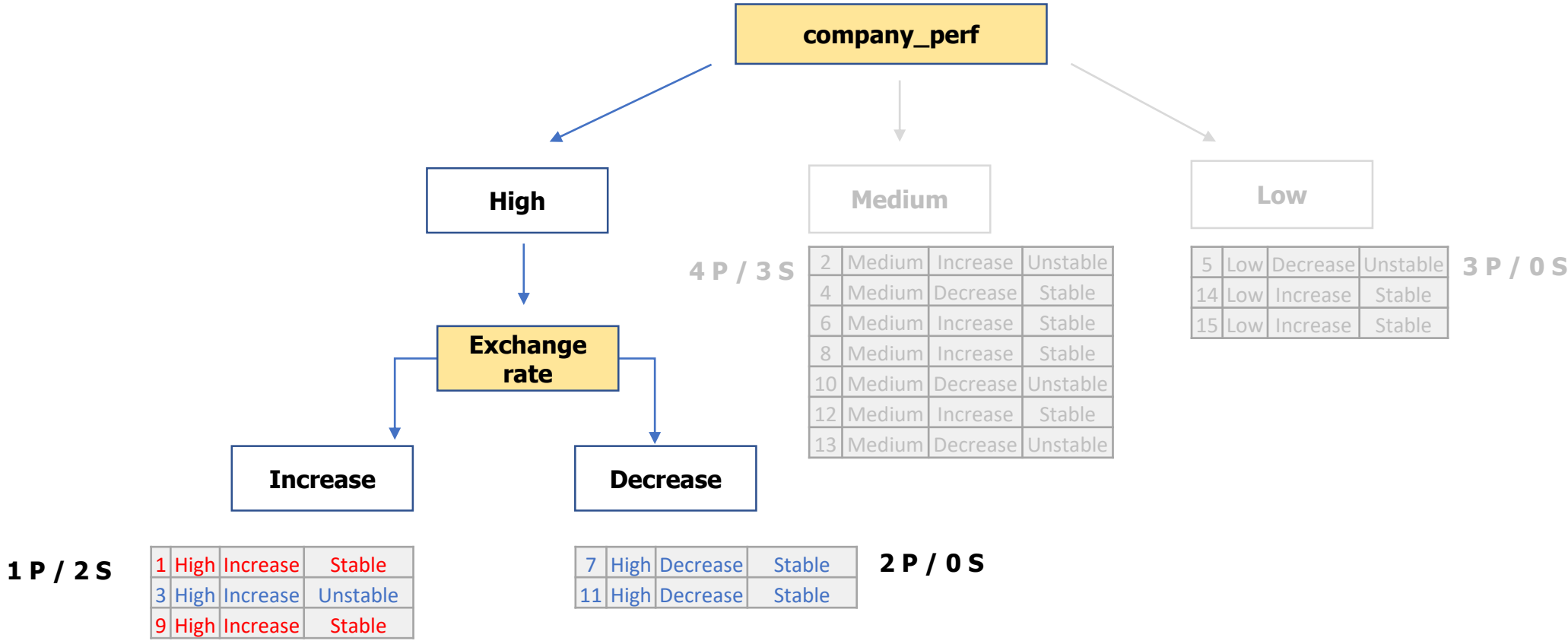
# Predict whether someone will buy or sell stocks

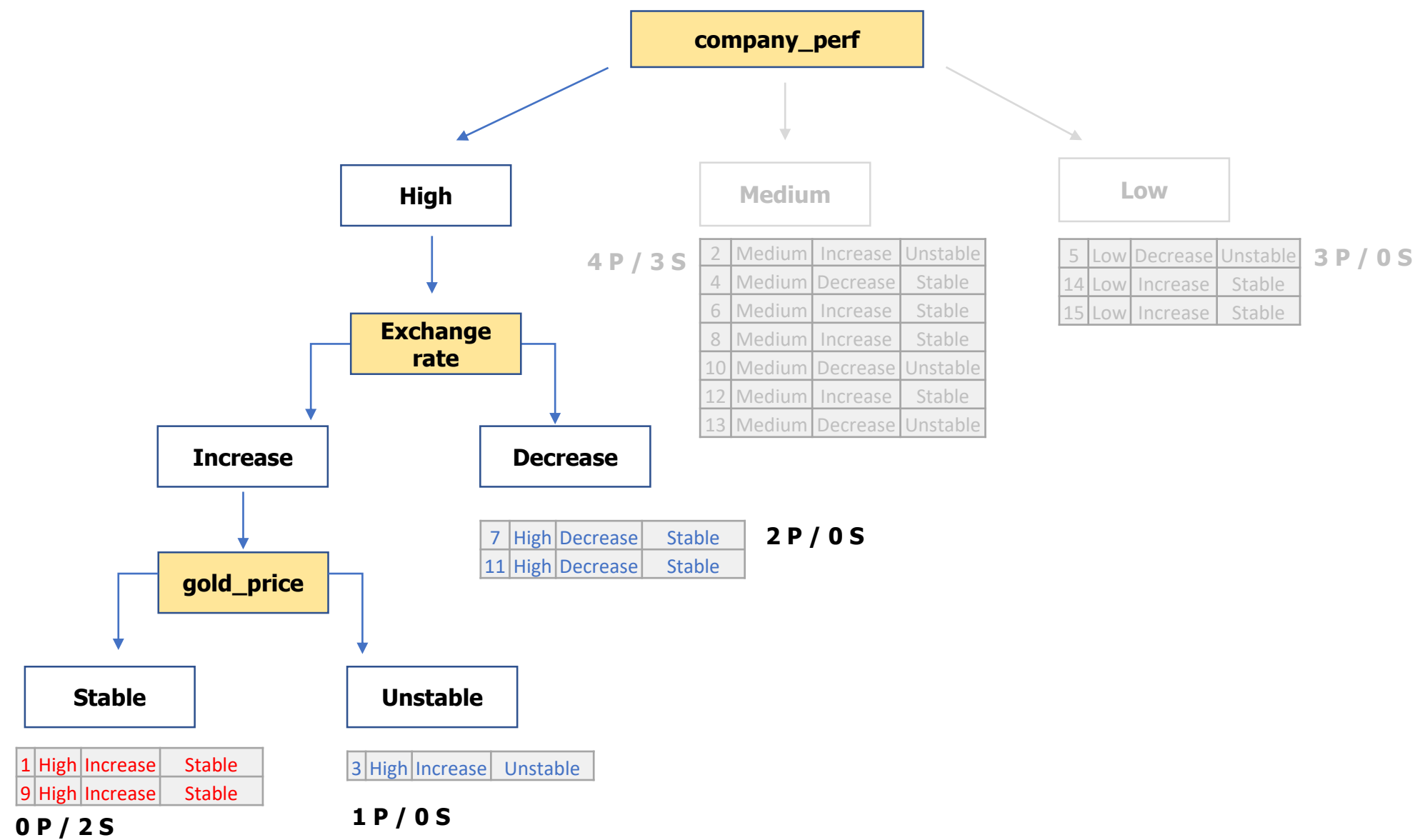
Day	company_perf	exchange_rate	gold_price	Action
1	High	Increase	Stable	Sale
2	Medium	Increase	Unstable	Purchase
3	High	Increase	Unstable	Purchase
4	Medium	Decrease	Stable	Purchase
5	Low	Decrease	Unstable	Purchase
6	Medium	Increase	Stable	Sale
7	High	Decrease	Stable	Purchase
8	Medium	Increase	Stable	Sale
9	High	Increase	Stable	Sale
10	Medium	Decrease	Unstable	Purchase
11	High	Decrease	Stable	Purchase
12	Medium	Increase	Stable	Sale
13	Medium	Decrease	Unstable	Purchase
14	Low	Increase	Stable	Purchase
15	Low	Increase	Stable	Purchase
16	Low	Decrease	Stable	????

**Company\_perf** High, Medium, Low  
**Exch\_rate** Increase, Decrease  
**Gold\_price** Stable, Unstable

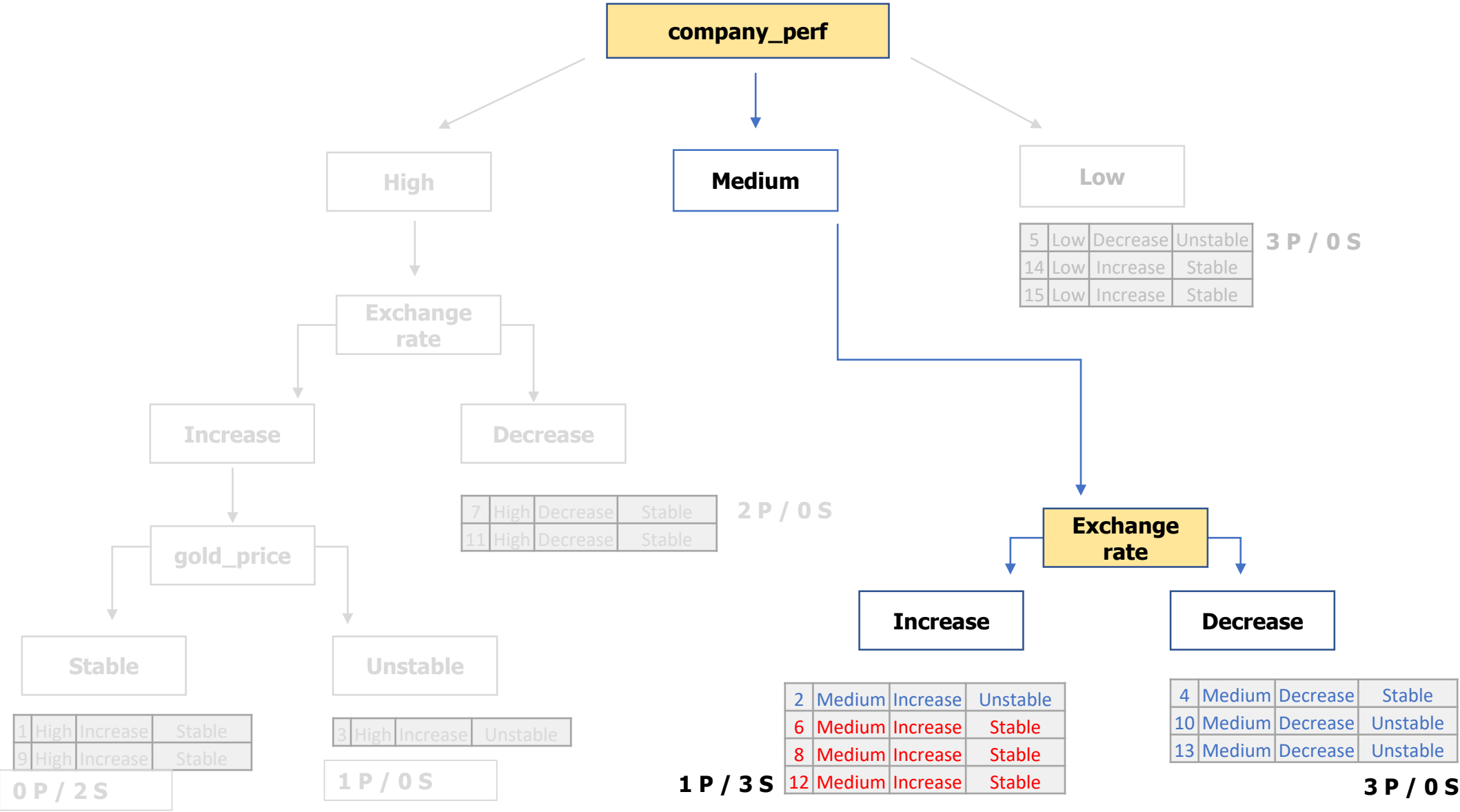
**Action** 9 Purchase / 6 Sale

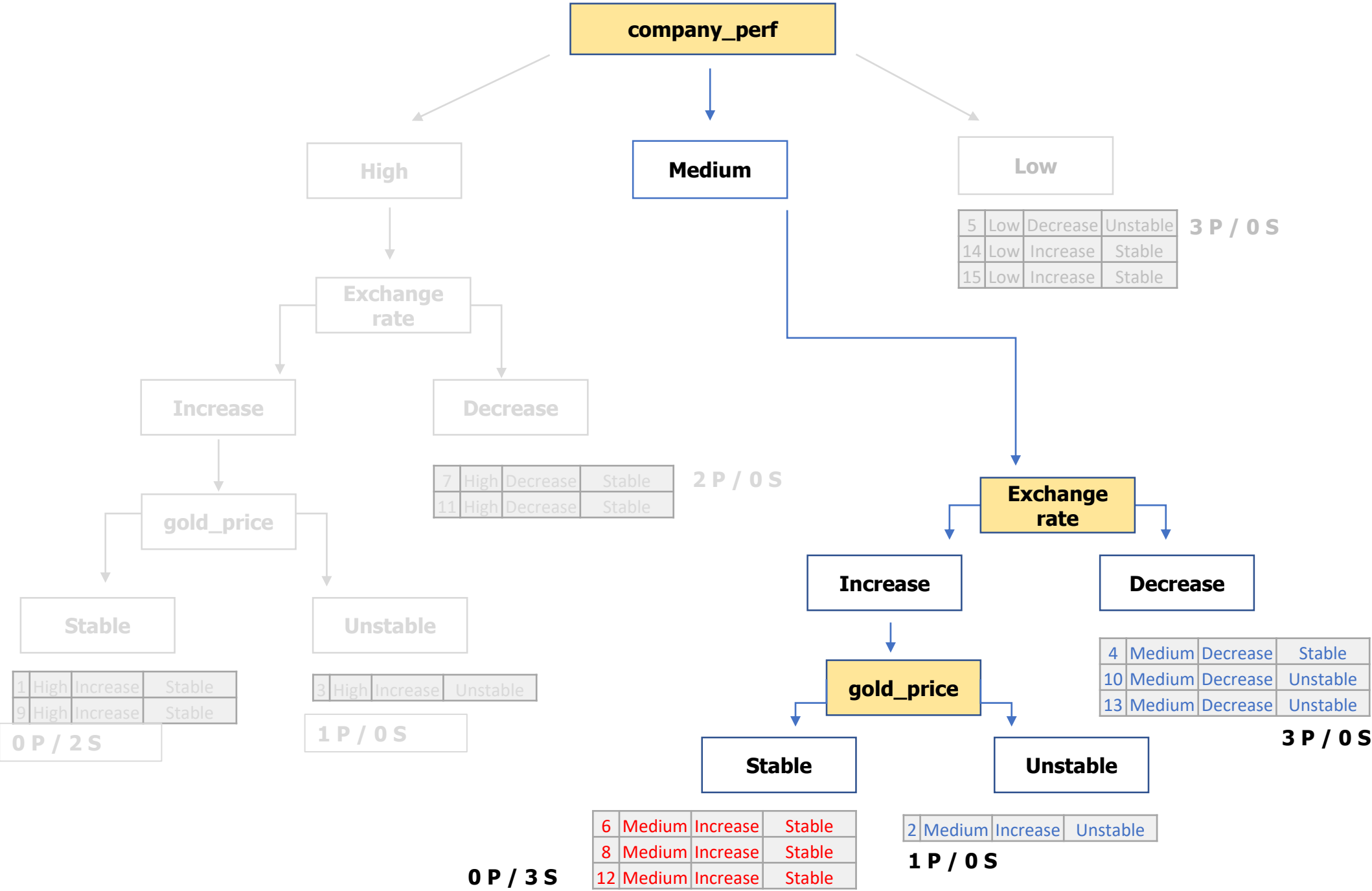




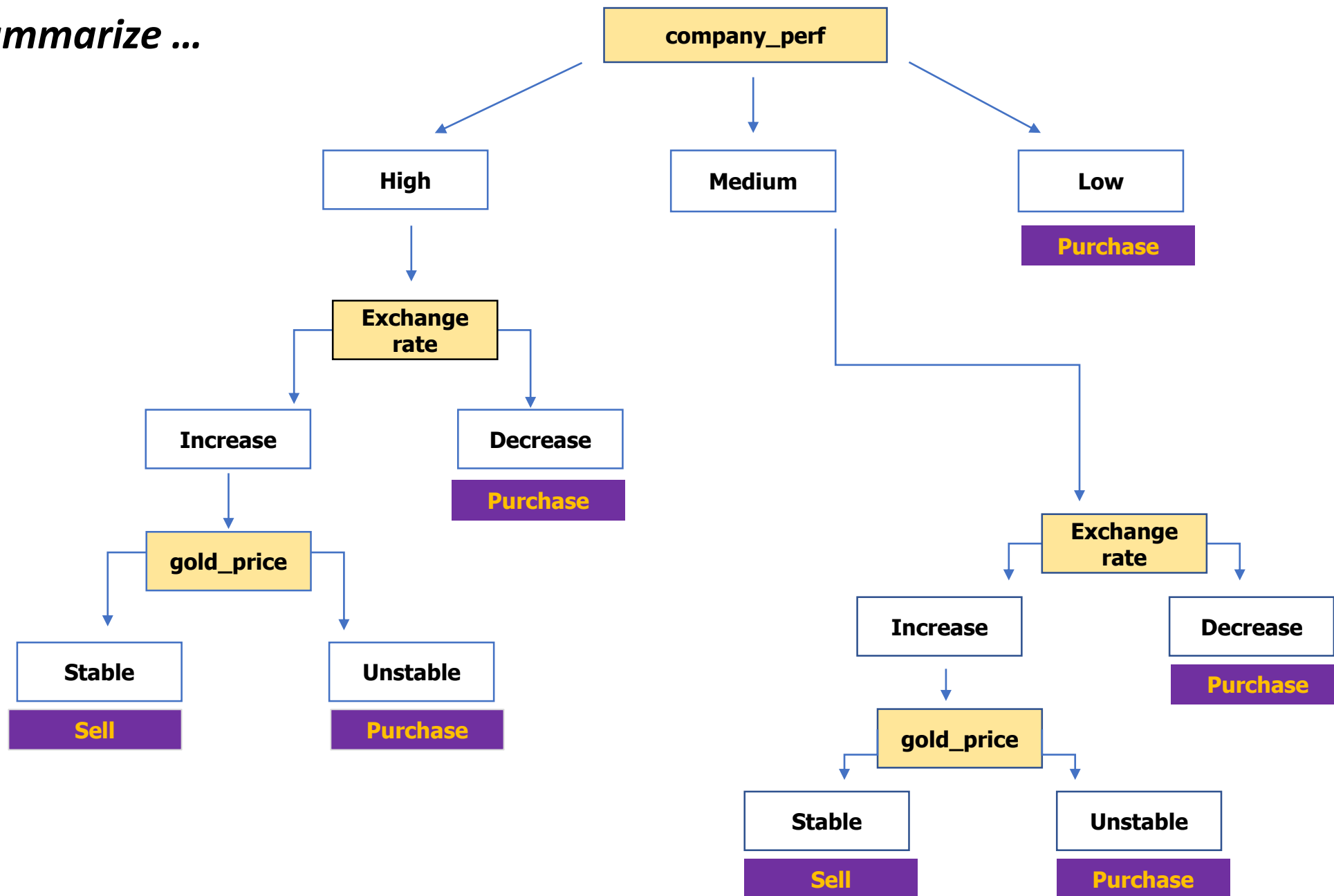


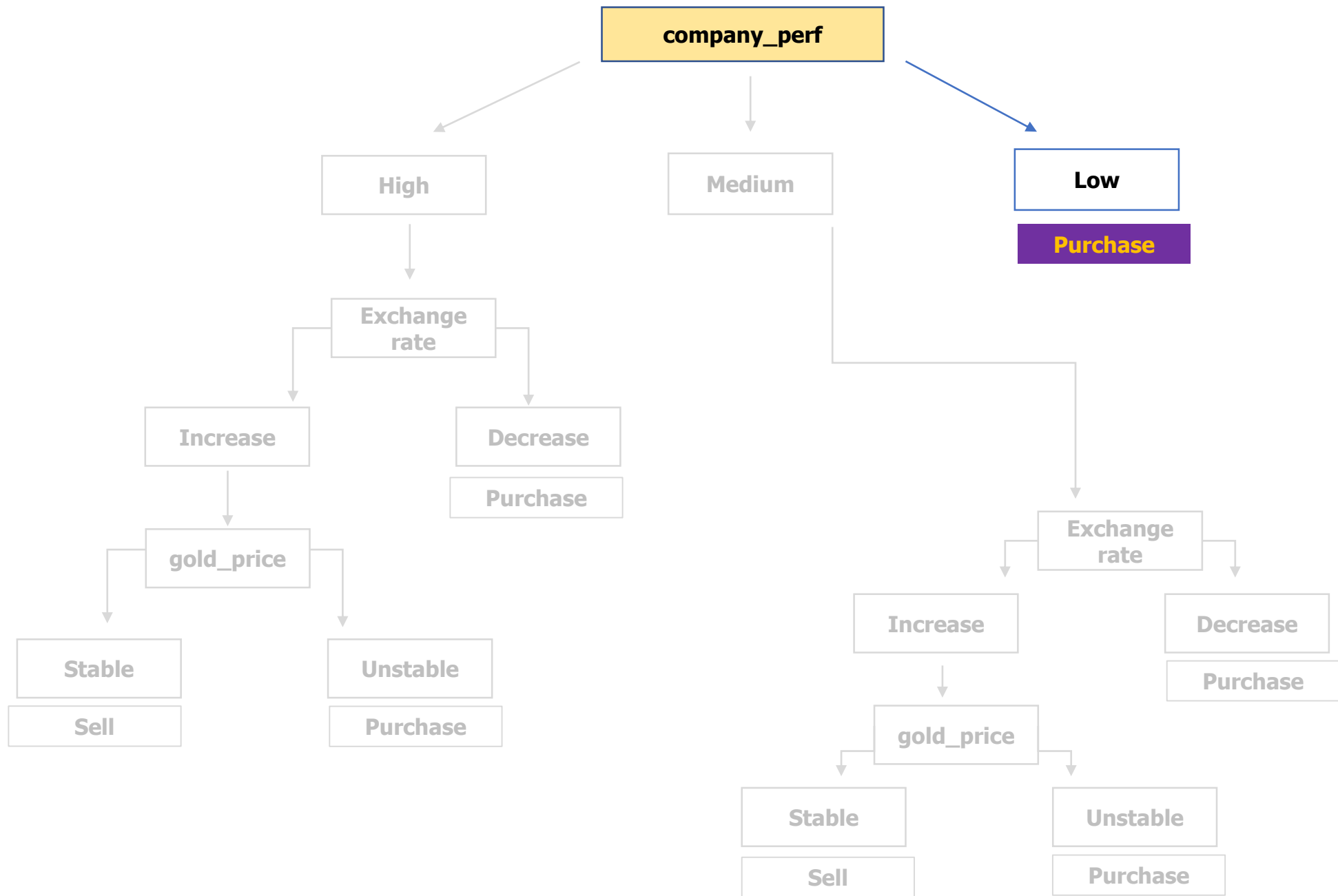






*To summarize ...*





Day	company_perf	exchange_rate	gold_price	Action
16	Low	Decrease	Stable	???

# Decision Trees - 2

- Conditions for split can be given (during the model building process). For example:
  - ☐ Don't split node if an attribute has same values
  - ☐ Don't split node if the number of children  $< 2$
  - ☐ Etc...

*rpart.control()* function in R lets you to specify the different conditions

- **Advantages**
  - ☐ Interpretable
  - ☐ Easy to understand
  - ☐ Scalable
  - ☐ Robust
- With more features, Trees can grow large and may become difficult to understand
- Smaller trees have better accuracy than larger trees
- Test dataset may become difficult to generalise (*tips example*)

# Split criteria

## Entropy

Total information held relating to the target variable (Binary) (IDE / C4.5)

More information, better will be the result

### Entropy (I) of the target variable

- ❑ Measures homogeneity of the sets
- ❑ Tells us how pure / impure a set is
- ❑ e.g. In a binary classification dataset, if **S** is the dataset having + and – classes, then Entropy (**I**nformation) is measured as:

$$E(S) = -p_{(+)} \log_2 p_{(+)} - p_{(-)} \log_2 p_{(-)}$$

where

$p_{(+)}$  = % of positive class

$p_{(-)}$  = % of negative class

- Interpretation of Entropy
  - ✓  $0 \leq I \leq 1$
  - ✓ Number of bits that is needed to identify if an item in the given dataset is + or –
  - ✓ For a pure subset, number of bits = 0
  - ✓ For a tie, number of bits = 1

### Information Gain

- ❑ Significant variable to split is determined by **Information Gain**
- ❑ Measure that determines how well a given attribute splits the dataset
- ❑ This measure is used at every step to determine the next best attribute
- ❑ Information (I) is needed to classify an object

$$\text{Gain}(S, A) = E(S) - \sum \left[ \left( \frac{S_a}{S} \right) * E(S_a) \right] \quad \text{(residual)}$$

where

$E(S)$  = Entropy calculation

$S_a$  = Count of attribute value **a**

$S$  = Total count of dataset of attribute **A**

$E(S_a)$  = Entropy of Attribute value **a**

- ❑ **Maximum**(Gain(A) ) → Best Attribute

# Which attribute to split on?

10 P / 5 S

company_perf	Purchase	Sale
High	3	2
Low	3	0
Medium	4	3

company_perf	Purchase (+)	Sale (-)	PS_Total	$p_{(+)}$	$\log_2 p_{(+)}$	$-p_{(+)} \log_2 p_{(+)}$	$p_{(-)}$	$\log_2 p_{(-)}$	$p_{(-)} \log_2 p_{(-)}$	E
High	3	2	5	0.600	-0.737	0.442	0.400	-1.322	-0.529	0.971
Low	3	0	3	1.000	0.000	0.000	0.000	0.000	0.000	0.000
Medium	4	3	7	0.571	-0.807	0.461	0.429	-1.222	-0.524	0.985
Set (Total)	10	5	15	0.667	-0.585	0.390	0.333	-1.585	-0.528	0.918
E(S)	S	$S_{(v=high)}$	$E_{(v=high)}$	$(S_a/S) * E(S_a)$	$S_{(v=medium)}$	$E_{(v=medium)}$	$(S_a/S) * E(S_a)$	Information Gain	Gain (company_perf)	
0.918	15	5	0.971	0.324	7	0.985	0.460	0.783	0.452	

10 P / 5 S

exchange_rate	Purchase	Sale
Decrease	6	0
Increase	4	5

exchange_rate	Purchase (+)	Sale (-)	PS_Total	$p_{(+)}$	$\log_2 p_{(+)}$	$-p_{(+)} \log_2 p_{(+)}$	$p_{(-)}$	$\log_2 p_{(-)}$	$p_{(-)} \log_2 p_{(-)}$	E
Decrease	6	0	6	1.000	0.000	0.000	0.000	0.000	0.000	0.000
Increase	4	5	9	0.444	-1.170	0.520	0.556	-0.848	-0.471	0.991
Set (Total)	10	5	15	0.667	-0.585	0.390	0.333	-1.585	-0.528	0.918
E(S)	S	$S_{(v=decrease)}$	$E_{(v=decrease)}$	$(S_a/S) * E(S_a)$	$S_{(v=increase)}$	$E_{(v=increase)}$	$(S_a/S) * E(S_a)$	Gain(S, exchange_rate)		
0.918	15	6	0.000	0.000	9	0.991	0.595	0.595		

10 P / 5 S

gold_price	Purchase	Sale
Stable	5	5
Unstable	5	0

gold_price	Purchase (+)	Sale (-)	PS_Total	$p_{(+)}$	$\log_2 p_{(+)}$	$-p_{(+)} \log_2 p_{(+)}$	$p_{(-)}$	$\log_2 p_{(-)}$	$p_{(-)} \log_2 p_{(-)}$	E
stable	5	5	10	0.500	-1.000	0.500	0.500	-1.000	-0.500	1.000
unstable	5	0	5	1.000	0.000	0.000	0.000	0.000	0.000	0.000
Set (Total)	10	5	15	0.667	-0.585	0.390	0.333	-1.585	-0.528	0.918
E(S)	S	$S_{(v=stable)}$	$E_{(v=stable)}$	$(S_a/S) * E(S_a)$	$S_{(v=unstable)}$	$E_{(v=unstable)}$	$(S_a/S) * E(S_a)$	Gain(S, company_perf)		
0.918	15	10	1.000	0.667	5	0.000	0.000	0.667		



Pure

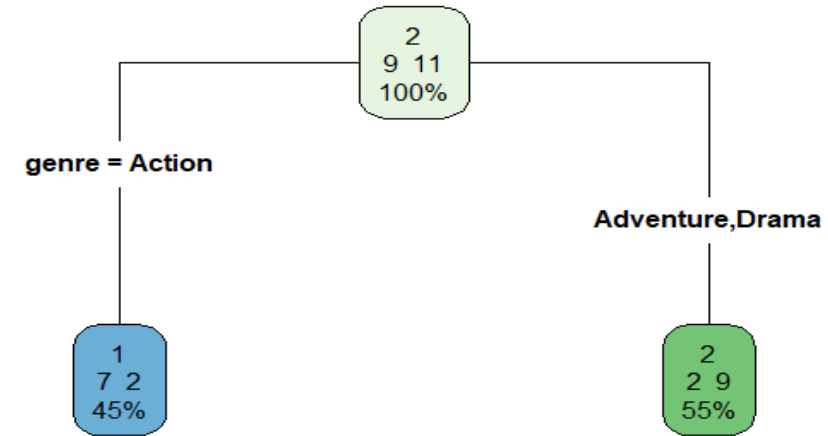


Impure

[Click here to view the detailed calculations](#)

# Exercise

#	creative_type	genre	rating	cat
1	Science Fiction	Action	PG-13	1
2	Fantasy	Adventure	PG-13	1
3	Fantasy	Adventure	PG	2
4	Fantasy	Drama	PG-13	2
5	Fantasy	Drama	PG-13	2
6	Science Fiction	Action	PG-13	1
7	Super Hero	Action	PG	1
8	Super Hero	Action	PG	1
9	Super Hero	Action	PG-13	2
10	Super Hero	Drama	R	2
11	Super Hero	Drama	PG-13	2
12	Science Fiction	Drama	PG-13	2
13	Science Fiction	Drama	R	2
14	Science Fiction	Action	PG	2
15	Science Fiction	Action	R	1
16	Fantasy	Action	R	1
17	Fantasy	Action	R	1
18	Fantasy	Adventure	R	1
19	Fantasy	Adventure	PG	2
20	Fantasy	Adventure	PG-13	2

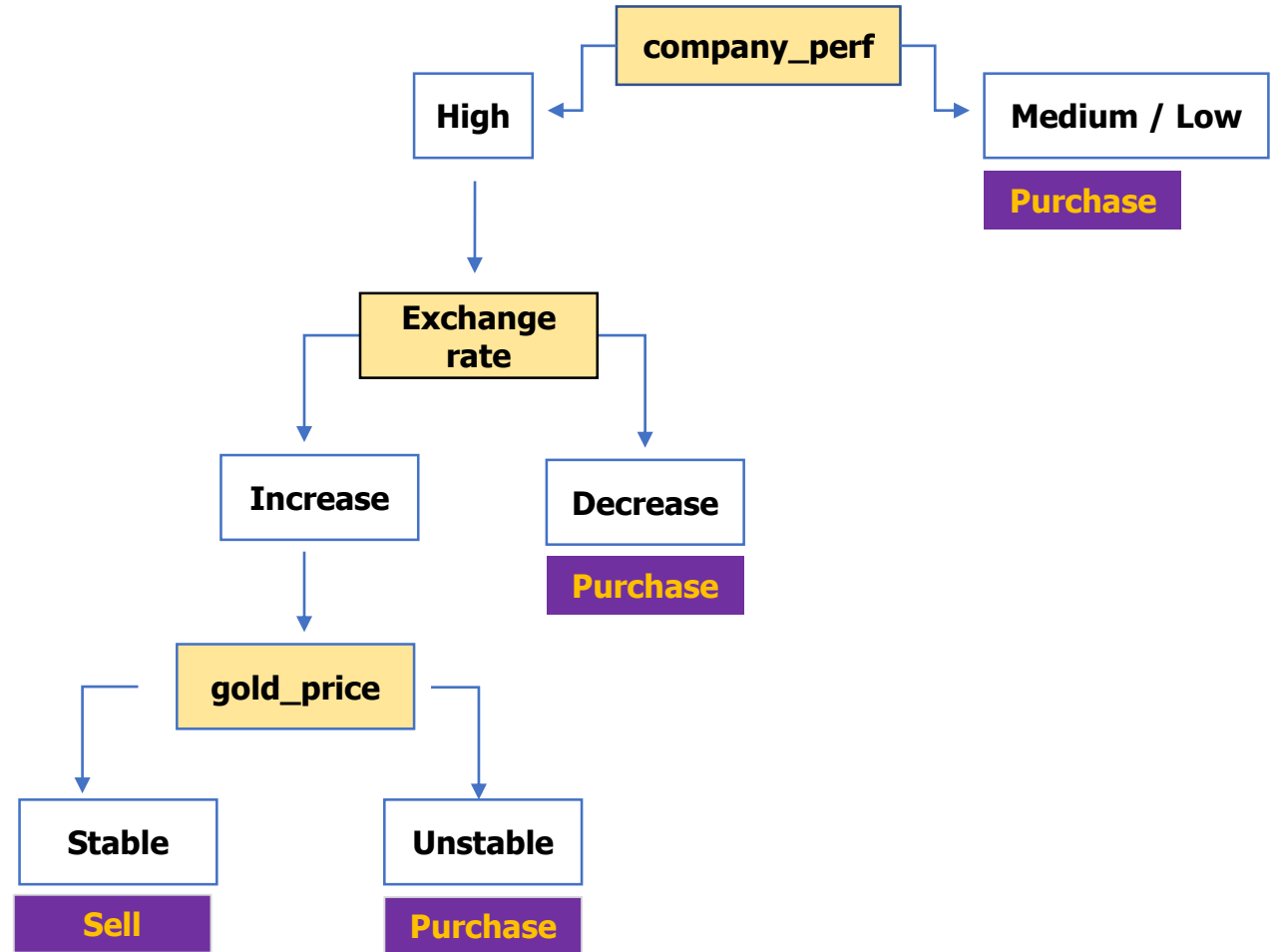


Find the first best attribute to split



# Pruning the Decision Tree

- Pruning is a technique to reduce the size of the Decision Tree by eliminating certain sections of the tree that provide little information to classify instances
- Pruning a tree is done to prevent overfitting; thereby improving accuracy
- Select the tree size that minimises the cross-validated error (*R has in-built function for this*)
- Pruning the Decision Tree
  - ❑ Pre-pruning : Chi-square test
  - ❑ Post-pruning : Pruning techniques to reduce the tree size (recommended)



- Pruning is done using a technique known as “Complexity Parameter”
- Post pruning, perform prediction with the pruned tree
- Compare the results with the pre-pruned model to check effectiveness
- `Plotcp()` provides a graphical representation of the cross-validation error summary

