# Árboles de decisión y ensambladores
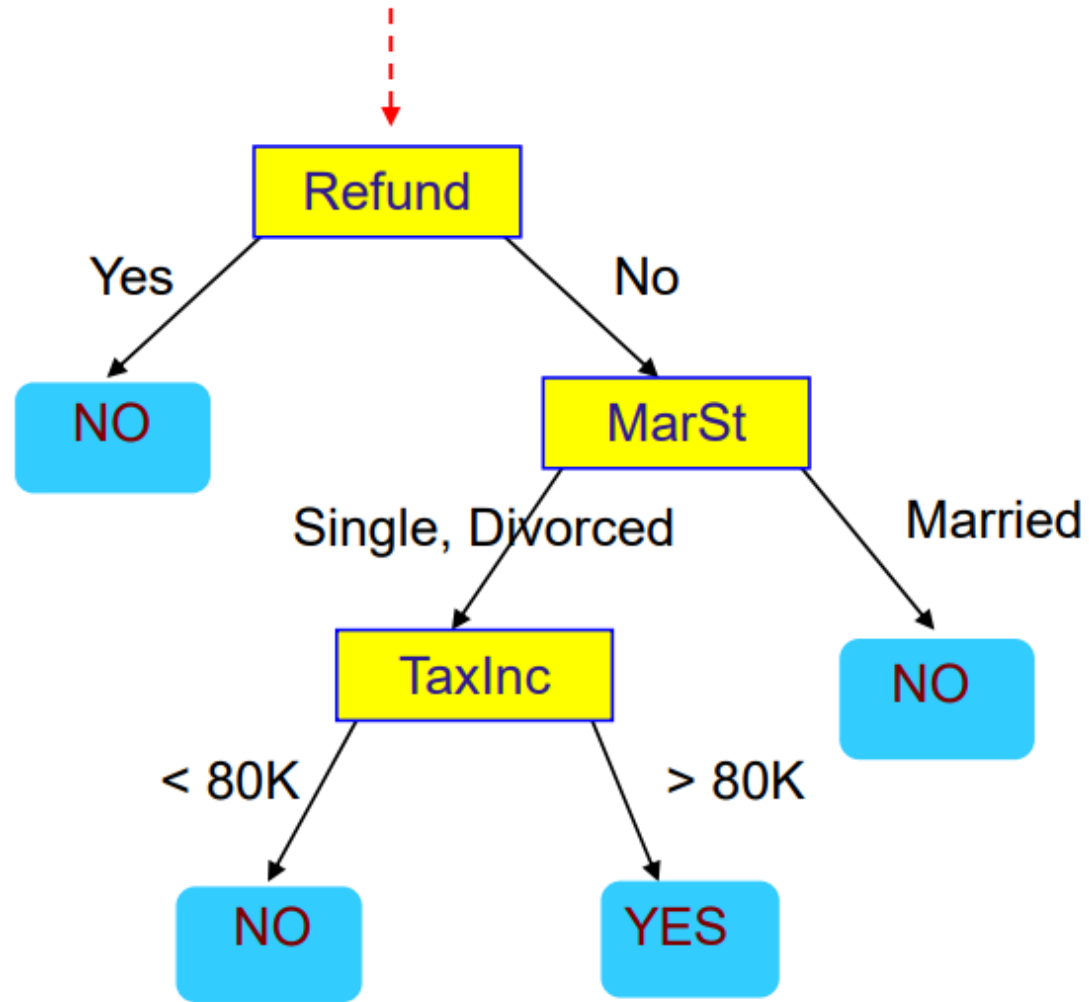
# Classification trees

# How to use



| Refund | Marital Status | Taxable Income | Cheat |
|--------|----------------|----------------|-------|
| No     | Married        | 80K            | ?     |

# Aplicando el modelo a los datos de evaluación



| Refund | Marital Status | Taxable Income | Cheat |
|--------|---------------|----------------|-------|
| No | Married | 80K | ? |

Test Data

Assign Cheat to "No"

# Impurity Criterion

## Gini Index

$$I_G = 1 - \sum_{j=1}^{c} p_j^2$$

$p_j$: proportion of the samples that belongs to class c for a particular node

## Entropy

$$I_H = - \sum_{j=1}^{c} p_j log_2(p_j)$$

$p_j$: proportion of the samples that belongs to class c for a particular node.

*This is the the definition of entropy for all non-empty classes ($p \neq 0$). The entropy is 0 if all samples at a node belong to the same class.
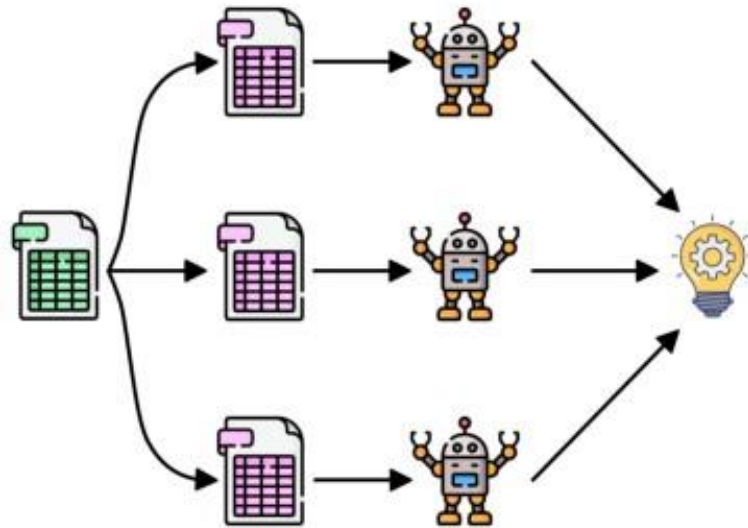
$$E(S) = \sum_{i=1}^{c} - p_i \log_2 p_i$$

| Play Golf | |
|-----------|----|
| Yes | No |
| 9 | 5 |

**Entropy(PlayGolf)** = Entropy (5,9)

= Entropy (0.36, 0.64)

= - (0.36 $\log_2$ 0.36) - (0.64 $\log_2$ 0.64)

= 0.94

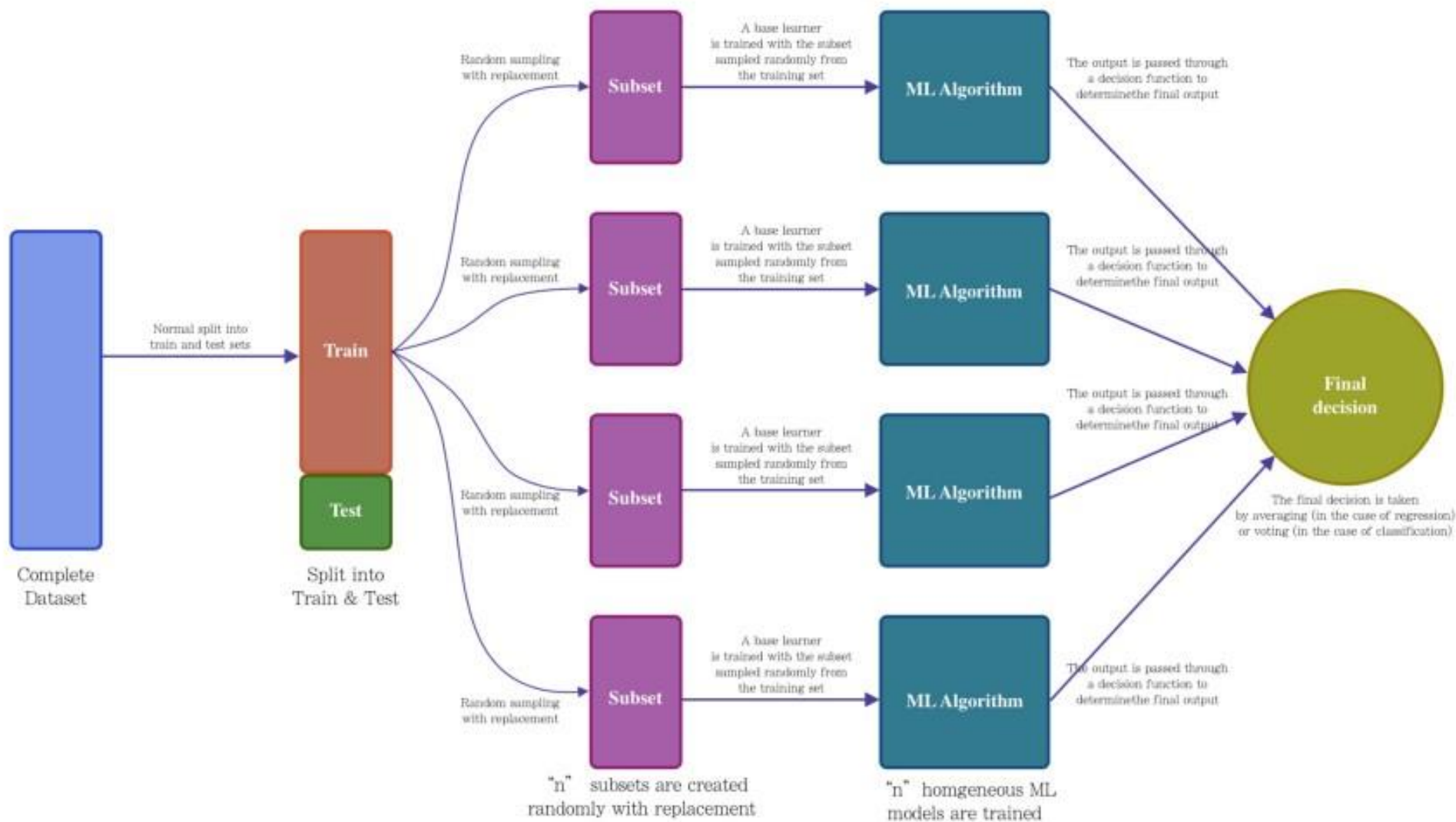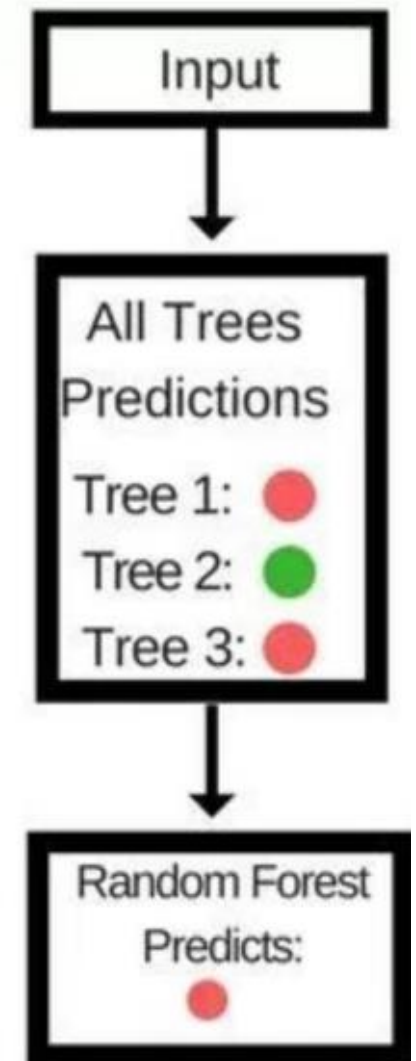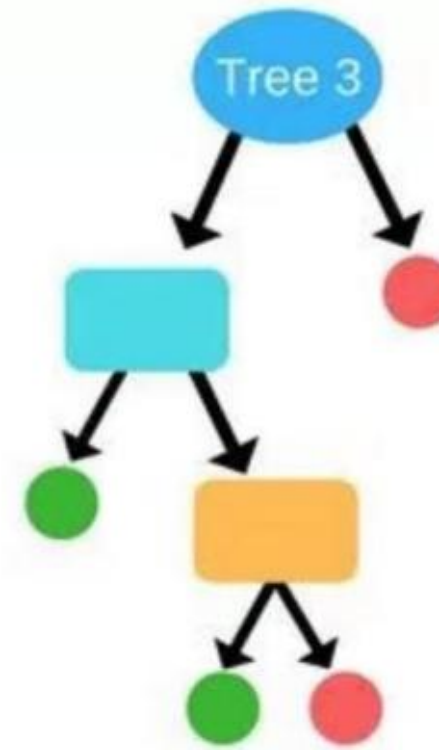# Bagging y boosting

# Bagging



Complete Dataset

Split into Train & Test

Normal split into train and test sets

Random sampling with replacement

A base learner is trained with the subset sampled randomly from the training set

The output is passed through a decision function to determine the final output

"n" subsets are created randomly with replacement

"n" homgeneous ML models are trained

Final decision

The final decision is taken by averaging (in the case of regression) or voting (in the case of classification)
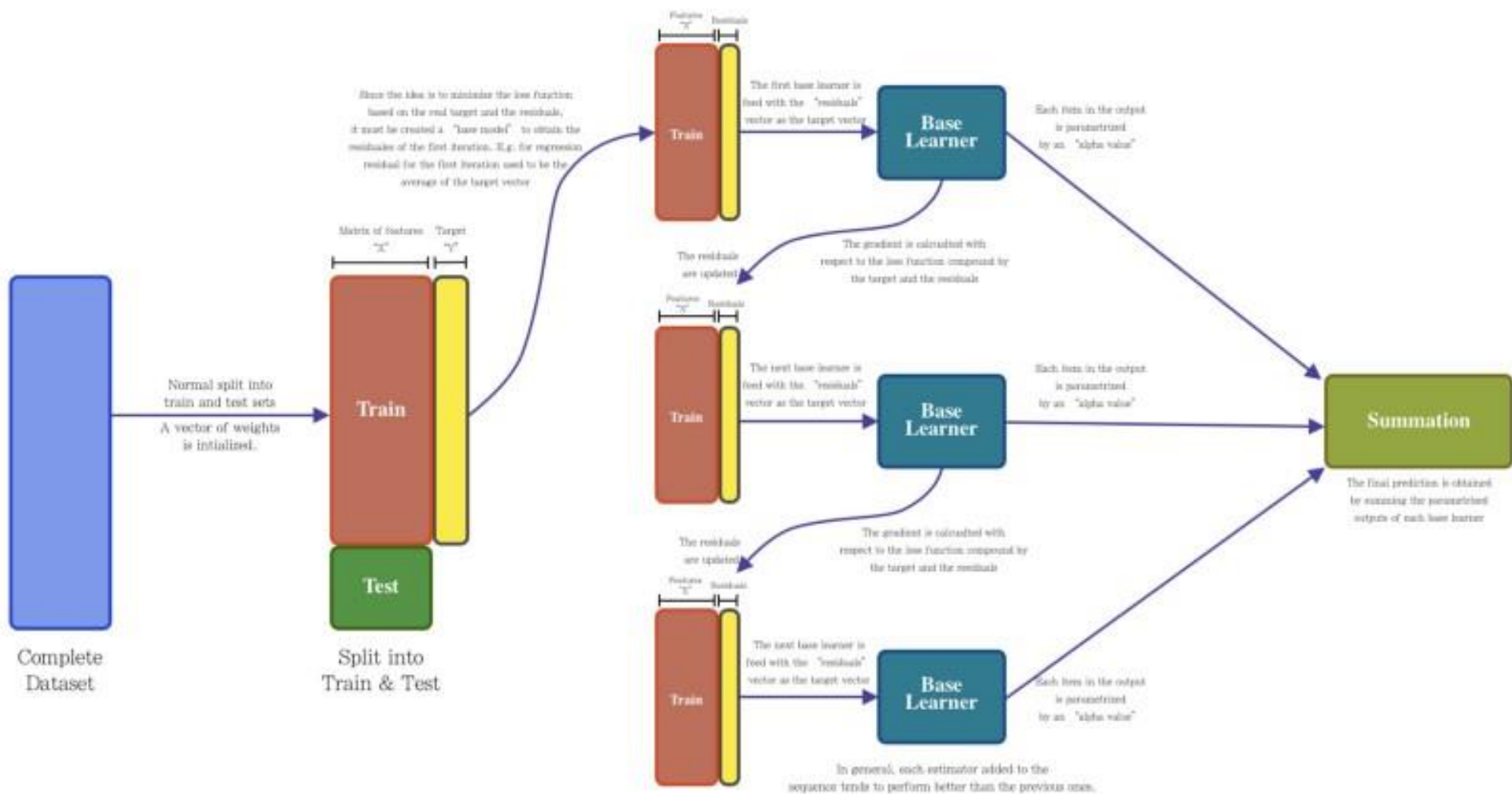
# Random Forest

# Boosting

# Random Forest vs Gradient Boosting

| Random Forest | Gradient Boosting |
|---|---|
| Easier to tune<br>Harder to overfit<br>Show very low variance<br>Easier to parallelize | Better accuracy with less trees if the data is noisy exhibits higher variance |

# Ensembling summary