

Interview Questions

1) What are some common hyperparameters of decision tree models, and how do they affect the model's performance?

①. Max depth:

- * Controls the maximum depth of trees

- * A deeper tree captures more patterns but can lead to overfitting.

②. Min Sample Split

- * The minimum number of samples requires to split an internal node.

- * Higher values prevent overfitting by ensuring split.

③. Min Sample Leaf:-

- * The minimum number of samples a leaf node must contain.

- * A smaller number reduces overfitting and increases generalization but may also lower accuracy.

④. Max Features:-

- * The number of features to consider when looking for the best split.

⑤. Criterion:-

- * Gini measures impurity & entropy

* entropy:- measure the information gain

* MSE (mean squared Error):- For regression tree

⑥. Max Leaf Nodes:-

* limits the number of leaf nodes. ①

* Helps prevent overfitting by controlling tree complexity

⑦. Min Impurity Decrease:-

* Prevents unnecessary split and controls overfitting. ②

Effect on Model Performance:-

① Overfitting vs underfitting

* Small, Large, large lead to underfitting
* Large (max depth), small (min sample split) leads to overfitting.

② Bias-variance Tradeoff

* High depth & small constraints reduce bias but increase variance.

* Regularization through constraints

③ Computational Efficiency:-

* Large trees are more computationally expensive.

* Limiting & max leaf nodes improve efficiency.

2) What is different between ~~En~~ label Encoding & one-hot-encoding?

1) Label-Encoding.

* Assigns a unique integer category.

Ex:-

categories: ['Saran', 'Harri', 'Deepak']

Label Encoded: {'Saran': 0, 'Harri': 1}

{ 'Deepak': 2 }

* Pros:-

* Simple & memory-efficient

* Cons:-

* Imposes an ordinal relationship between categories, which may not be meaningful.

* use case:-

* Suitable for ordinal categorical variable, where order matters (low to high)

2) One-Hot-Encoding:-

* Convert each category into a separate binary column (0 or 1)

Ex:-

categories: ['Red', 'Green', 'Blue']

One-Hot-Encoded: 1 0 0

Red Green Blue

1 0 0

propagator register output

0 0 1

['dog', 'cat', 'bird']

* Pros:-

* No ordinary relationship assumption.

* Cons:-

* Increases dimensionality,

leading to the curse of dimensionality for large datasets.

* Use Case:-

* Suitable for nominal categorical variable.

at what? reform reform, reform

One-Hot-Encoded: 1 0 0

['dog', 'cat', 'bird']