

① Explain Bagging and Boosting methods.

How is it different from each other.

Bagging vs Boosting understanding
so different, no basis is there*

Bagging & Boosting are two popular ensemble learning techniques used to improve the performance and accuracy of machine learning model. While both methods combine multiple weaker learners to create a strong predictive

1). Bagging:

Bagging aim reduce variance by training multiple models independently on different subsets of the datasets and then averaging their predictions.

Key features: *

* uses Bootstrap to create different subsets of data and decisions

* Trains multiple independent models in parallel and a majority

Ex: ~~Random forest~~ has pruned trees

- * Random forest is an extension of ~~pruning~~ applied to decision trees or pruned subtrees of ~~pruned~~ trees.
- * Each tree is trained on a different bootstrap sample.

Note: Out of bagged & pruned

2) Boosting involves parallel adaboost

In boosting rounds ~~it~~ assign weight to reduce bias by training a model learns from the mistakes of the previous ones. Note: A sign of overfitting is when a model is able to learn from the training data but fails to perform well on test data.

Key features:-

- * Model's are trained sequentially each correcting the mistakes of its predecessors. Note: Out of bagged
- * Can't overfit note it is fully interleaved.

~~It~~ it's adding the weights of correctly classified points to the loss function.

1) AdaBoost:-

- * Assign higher weight to misclassified point.

2) Gradient Boosting:-

- * Models learn sequentially by minimizing a loss function.

② Explain how to handle imbalance in the data. with diagram. Label

Handling imbalanced Data in ML

Imbalanced data occurs when one class significantly outnumbers another in a classification problem.

1) Data - Entry Techniques

These techniques deal with imbalanced class distribution.

A) OverSampling

* SMOTE generates synthetic samples for the minority class by interpolating existing instances.

* ADASYN is similar to SMOTE but focuses more on harder-to-learn minority instances.

* Random oversampling - Duplicates existing minority samples.

B) UnderSampling

* Random underSampling

Remove random samples from the majority class.

* Nearmiss Algorithm :- at least one misclassified
Select majority class samples

2) Algorithm level Techniques:-
Instead of modifying data,
these methods fit a function such
function. making minority as a
function.

A) Cost - Sensitive Learning:

* Assign better misclassification
costs to the minority. makes the
model more sensitive to errors.

from Sklearn. ensemble import
BalancedRandomForestClassifier

RFC: Random Forest Classifier (\rightarrow Ensemble
and Tree of Maximum weight:
most of reward no evatt. aspect
"balanced"

B) Ensemble Techniques.

* Balanced Random Forest:-
uses bootstrapping, picking
equals class representation

* Easy Ensemble:-
combine bootstrapping, under sampling &
over sampling

3) Evaluation metrics for Imbalanced data.

- * Precision & Recall

- * F1-Score

- * ROC-AUC

- * PR-AUC

from sklearn.metrics import classification_report.

classification_report(y_Pred, y_test)

4) Alternative Approaches:-

- * Anomaly Detection:-

The minority very rare,
treat anomaly detection problem.

Choosing the right methods.

- * Small dataset - SMOTE -

sensitive learning

- * Large dataset - under sampling +
ensemble method.

- * Extreme imbalance :- consider anomaly
detection techniques.