Question 1: What is Ensemble Learning in machine learning? Explain the key idea behind it.

Ans:-Ensemble Learning in machine learning is a technique where multiple models (often called base learners or weak learners) are combined to solve a problem and improve overall performance compared to using a single model.



Key Idea Behind Ensemble Learning:

The key idea is that a group of models, when combined, can produce more accurate and reliable predictions than any individual model alone.

- Different models may make different errors.
- By combining their predictions, the ensemble can cancel out individual errors, leading to better generalization.
- It's based on the principle that "the wisdom of the crowd" often outperforms a single expert.

Question 2: What is the difference between Bagging and Boosting?

Aspect	Bagging (Bootstrap Aggregating)	Boosting
Purpose	Reduce variance	Reduce bias and variance
Model Training	Models are trained independently and in parallel	Models are trained sequentially , each learning
Data Sampling	Each model gets a random subset of data (with replacement)	Each new model focuses on misclassified example.
Weighting	All models have equal weight in final prediction	Models are weighted based on their accuracy
Overfitting	Less prone to overfitting	More prone to overfitting (but better generaliza
Example Algorithm	Random Forest	AdaBoost, Gradient Boosting, XGBoost
Prediction Output	Majority vote (classification) or average (regression)	Weighted vote or sum of model predictions

Question 3: What is bootstrap sampling and what role does it play in Bagging methods like Random Forest?

Ans:- Bootstrap sampling is a method of creating new datasets by randomly selecting data points with replacement from the original dataset. This means that:

Each bootstrap sample is the same size as the original dataset.

Some data points may appear multiple times, and others may not appear at all.

Marks: How Bootstrap Sampling Works:

Given a dataset of size N, bootstrap sampling involves:

Randomly selecting N data points, with replacement, from the original dataset.

Creating a new dataset (called a "bootstrap sample") of size N.

Repeating this process to generate multiple bootstrap samples.

® Role in Bagging (e.g., Random Forest):

In Bagging (Bootstrap Aggregating), and specifically in Random Forest, bootstrap sampling plays a crucial role:

- Training Diversity:
 - Each decision tree in the forest is trained on a different bootstrap sample.
 - This introduces variation among the trees, which helps reduce overfitting.
- Model Independence:
 - Since each tree sees a different subset of the data, their predictions are less correlated.
 - Combining their outputs (via majority vote or averaging) leads to better generalization.
- ✓ Out-of-Bag (OOB) Error Estimation:
 - On average, ~63% of the original data appears in each bootstrap sample.
 - The remaining ~37% can be used as validation data to estimate model accuracy without using a separate test set.

Question 4: What are Out-of-Bag (OOB) samples and how is OOB score used to evaluate ensemble models?

Ans:- In ensemble methods like Bagging (e.g., Random Forest), Out-of-Bag (OOB) samples refer to the data points not included in a given bootstrap sample.

③ What is the OOB Score?

The OOB score is an internal cross-validation accuracy estimate that uses the OOB samples.

How It's Used:

- 1.Each model is trained on its bootstrap sample.
- 2.Its OOB samples are used as a validation set.
- 3.After all models are trained:
 - Each data point may be OOB for several models.
 - These models make predictions on the data point.
 - The final prediction is typically made by majority vote (classification) or average (regression).
- 4.The OOB score is the accuracy (or R² for regression) of these predictions compared to the true labels.

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.datasets import load_iris
# Load dataset
iris = load_iris()
X, y = iris.data, iris.target
# Train Random Forest with OOB scoring enabled
rf = RandomForestClassifier(n_estimators=100, oob_score=True, bootstrap=True, random
rf.fit(X, y)
# OOB score
print("00B Score:", rf.oob_score_)
```

→ 00B Score: 0.9533333333333334

Question 5: Compare feature importance analysis in a single Decision Tree vs. a Random Forest.

Ans:- Feature importance analysis helps us understand which input features are most influential in making predictions.

In a Single Decision Tree:

- How It Works:

- Feature importance is calculated based on how much each feature reduces impurity (like Gini Impurity or Entropy) at each split.
- The more a feature is used to split the data and reduce impurity, the more important it is considered.

Pros:

- Simple to compute and interpret.
- Reflects the direct structure of the tree.

Cons:

- High variance: A single tree can change drastically with small changes in data.
- May overemphasize certain features due to overfitting

In a Random Forest (Ensemble of Trees):

How It Works:

- Feature importance is averaged across all the trees in the forest.
- Each tree calculates its own feature importance, and then the results are aggregated (typically normalized).

Pros:

- More stable and robust than a single tree.
- Reduces the risk of overestimating the importance of irrelevant features.
- Better at generalization.

Cons:

- Harder to interpret than a single tree (less transparency).
- Can be biased toward features with more categories or higher cardinality.

Question 6: Write a Python program to: • Load the Breast Cancer dataset using sklearn.datasets.load_breast_cancer() • Train a Random Forest Classifier • Print the top 5 most important features based on feature importance scores

Ans:- **%** Python Code:

```
from sklearn.datasets import load_breast_cancer
from sklearn.ensemble import RandomForestClassifier
import pandas as pd
# Step 1: Load the Breast Cancer dataset
data = load_breast_cancer()
X = data.data
y = data.target
feature_names = data.feature_names
# Step 2: Train a Random Forest Classifier
rf = RandomForestClassifier(n_estimators=100, random_state=42)
rf.fit(X, y)
# Step 3: Get feature importance scores
importances = rf.feature_importances_
# Step 4: Create a DataFrame for better readability
feature_importance_df = pd.DataFrame({
    'Feature': feature_names,
    'Importance': importances
})
# Step 5: Sort and print the top 5 most important features
top_5_features = feature_importance_df.sort_values(by='Importance', ascending=False)
print(" Top 5 Most Important Features:")
print(top_5_features)
```

```
23 worst area 0.139357
27 worst concave points 0.132225
7 mean concave points 0.107046
20 worst radius 0.082848
22 worst perimeter 0.080850
```

Question 7: Write a Python program to:

Train a Bagging Classifier using Decision Trees on the Iris dataset
 Evaluate its accuracy and compare with a single Decision Tree

Ans:- V Python Code:

```
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import BaggingClassifier
from sklearn.metrics import accuracy_score
# Load the Iris dataset
iris = load iris()
X = iris.data
y = iris.target
# Split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_stat
# 1. Train a single Decision Tree Classifier
dt_classifier = DecisionTreeClassifier(random_state=42)
dt_classifier.fit(X_train, y_train)
# Make predictions and evaluate accuracy for the single Decision Tree
y_pred_dt = dt_classifier.predict(X_test)
accuracy_dt = accuracy_score(y_test, y_pred_dt)
print(f"Accuracy of single Decision Tree: {accuracy_dt:.4f}")
# 2. Train a Bagging Classifier using Decision Trees
# The base_estimator is set to DecisionTreeClassifier
# n_estimators determines the number of base estimators (decision trees)
bagging_classifier = BaggingClassifier(
    estimator=DecisionTreeClassifier(random_state=42),
    n_estimators=10, # Number of decision trees in the ensemble
    random_state=42
)
bagging_classifier.fit(X_train, y_train)
# Make predictions and evaluate accuracy for the Bagging Classifier
y_pred_bagging = bagging_classifier.predict(X_test)
accuracy_bagging = accuracy_score(y_test, y_pred_bagging)
print(f"Accuracy of Bagging Classifier: {accuracy_bagging:.4f}")
```

```
# Compare the accuracies
print("\nComparison:")
if accuracy_bagging > accuracy_dt:
    print("Bagging Classifier performed better than a single Decision Tree.")
elif accuracy_bagging < accuracy_dt:
    print("Single Decision Tree performed better than Bagging Classifier.")
else:
    print("Both models achieved the same accuracy.")

Accuracy of single Decision Tree: 1.0000
Accuracy of Bagging Classifier: 1.0000

Comparison:
Both models achieved the same accuracy.
```

Question 8: Write a Python program to: ● Train a Random Forest Classifier ● Tune hyperparameters max_depth and n_estimators using GridSearchCV ● Print the best parameters and final accuracy

```
from sklearn.datasets import load_iris
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import GridSearchCV, train_test_split
from sklearn.metrics import accuracy_score
# Step 1: Load dataset
data = load_iris()
X, y = data.data, data.target
# Step 2: Split into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_stat
# Step 3: Define the parameter grid
param_grid = {
    'n_estimators': [50, 100, 150],
    'max_depth': [3, 5, 7, None]
}
# Step 4: Create Random Forest Classifier
rf = RandomForestClassifier(random_state=42)
# Step 5: Use GridSearchCV for hyperparameter tuning
grid_search = GridSearchCV(estimator=rf, param_grid=param_grid, cv=5, n_jobs=-1, sco
grid_search.fit(X_train, y_train)
# Step 6: Get the best model and evaluate on test set
best_model = grid_search.best_estimator_
```

```
y_pred = best_model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)

# Step 7: Print the results
print(" Best Parameters:", grid_search.best_params_)
print(" Final Accuracy on Test Set: {:.2f}%".format(accuracy * 100))

** Best Parameters: {'max_depth': 3, 'n_estimators': 150}

** Final Accuracy on Test Set: 100.00%
```

Question 9: Write a Python program to: • Train a Bagging Regressor and a Random Forest Regressor on the California Housing dataset • Compare their Mean Squared Errors (MSE)

```
from sklearn.datasets import fetch_california_housing
from sklearn.model_selection import train_test_split
from sklearn.ensemble import BaggingRegressor, RandomForestRegressor
from sklearn.metrics import mean_squared_error
# Load the California Housing dataset
housing = fetch_california_housing(as_frame=True)
X = housing.data
y = housing.target
# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_stat
# Train a Bagging Regressor
bagging_regressor = BaggingRegressor(random_state=42)
bagging_regressor.fit(X_train, y_train)
y_pred_bagging = bagging_regressor.predict(X_test)
mse_bagging = mean_squared_error(y_test, y_pred_bagging)
# Train a Random Forest Regressor
random_forest_regressor = RandomForestRegressor(random_state=42)
random_forest_regressor.fit(X_train, y_train)
y_pred_rf = random_forest_regressor.predict(X_test)
mse_rf = mean_squared_error(y_test, y_pred_rf)
# Compare their Mean Squared Errors
print(f"Mean Squared Error (Bagging Regressor): {mse_bagging:.4f}")
print(f"Mean Squared Error (Random Forest Regressor): {mse_rf:.4f}")
```

Mean Squared Error (Bagging Regressor): 0.2824
Mean Squared Error (Random Forest Regressor): 0.2554

Question 10: You are working as a data scientist at a financial institution to predict loan default. You have access to customer demographic and transaction history data.

You decide to use ensemble techniques to increase model performance. Explain your step-bystep approach to:

 Choose between Bagging or Boosting ● Handle overfitting ● Select base models ● Evaluate performance using cross-validation ● Justify how ensemble learning improves decision-making in this real-world context

Ans:- Q Step 1: Choose Between Bagging or Boosting

- Bagging (e.g., Random Forest)
 - When to use: If the model suffers from high variance and the data has noise.
 - Pros: Reduces overfitting, works well with uncorrelated models, more stable.
- ✓ Boosting (e.g., XGBoost, LightGBM)
 - When to use: If the model suffers from high bias or needs to capture complex patterns.
 - Pros: Learns from previous errors, performs well on imbalanced or structured datasets.

choice:

- Start with Boosting (e.g., XGBoost or LightGBM), as it tends to outperform on structured tabular data like customer profiles and transaction history.
- Test Bagging as a benchmark (e.g., with Random Forest) to compare robustness.

- Use cross-validation to detect overfitting during training.
- Apply regularization in boosting models:
 - For example, tune learning_rate, max_depth, min_child_weight in XGBoost.
- Limit model complexity:
 - Use early stopping with validation sets.
- Feature selection:
 - Drop irrelevant or highly correlated features to avoid noise.
- Use ensemble averaging (as in Bagging) to reduce variance.

Step 3: Select Base Models

- For Bagging: Use Decision Trees as base learners (e.g., in Random Forest).
- For Boosting: Use shallow trees (e.g., depth 3-6) as base learners.

You may experiment with:

Logistic Regression (as a weak learner in some stacking models)

• Gradient Boosting with categorical encoding for high-cardinality features

Step 4: Evaluate Performance Using Cross-Validation

Use Stratified K-Fold Cross-Validation to preserve class distribution (especially if data is imbalanced).

Evaluate using:

AUC-ROC: Measures model's ability to distinguish defaulters vs non-defaulters.

Precision/Recall / F1 Score: Especially important if false negatives (missed defaults) are costly.

Confusion Matrix: Helps visualize types of errors.

Use GridSearchCV or RandomizedSearchCV for hyperparameter tuning.

Step 5: Justify Ensemble Learning for Real-World Decision-Making

Why Ensemble Learning?

Higher accuracy: Combines multiple models to produce better generalization.

Better risk prediction: Boosting techniques adapt to complex patterns in transaction histories.

Robustness to noise: Bagging methods like Random Forest handle noisy customer data effectively.

Improved trust: Feature importance from ensembles can highlight key risk factors (e.g., payment history, income level).

Reduced business risk: More accurate predictions reduce loan defaults, increasing profit and reducing losses.