

1. **What is the definition of a target function? In the sense of a real-life example, express the target function. How is a target function's fitness assessed?**

The target function, also known as the **objective function** or **the ground truth**, is the ideal function or relationship that a machine learning model aims to approximate or predict. It represents the true underlying mapping between input features and output values. In a real-life example, consider a spam email detection system. The target function would be the accurate classification of emails as either spam or not spam. The fitness of a target function is assessed by comparing the model's predictions with the actual known values. Metrics such as accuracy, precision, recall, or mean squared error can be used to measure the fitness of the target function.

2. **What are predictive models, and how do they work? What are descriptive types, and how do you use them? Examples of both types of models should be provided. Distinguish between these two forms of models.**

Predictive models aim to make predictions or forecasts based on input data. **They learn patterns and relationships in the data to generalize and make predictions on new**, unseen instances. Examples of predictive models include regression models (predicting house prices based on features) and classification models (predicting customer churn based on demographic data). Descriptive models, on the other hand, **focus on summarizing and understanding the data** or the underlying phenomena. They provide insights into the relationships, distributions, or patterns in the data. Examples of descriptive models include clustering algorithms (identifying customer segments based on purchasing behavior) and association rule mining (finding patterns in market basket data). The main difference is that predictive models aim to make predictions, while descriptive models aim to describe and summarize the data.

3. **Describe the method of assessing a classification model's efficiency in detail. Describe the various measurement parameters.**

Assessing the efficiency of a classification model involves evaluating its performance using various measurement parameters. Common evaluation metrics include:

- **Accuracy:** The proportion of correct predictions over the total number of predictions.
  - **Precision:** The proportion of true positive predictions over the total number of positive predictions.
  - **Recall (Sensitivity):** The proportion of true positive predictions over the total number of actual positive instances.
  - **F1-measure:** The harmonic mean of precision and recall, providing a balanced measure of both metrics.
  - **Confusion matrix:** A table representing the counts of true positive, true negative, false positive, and false negative predictions.
4. **In the sense of machine learning models, what is underfitting? What is the most common reason for underfitting?**

Underfitting occurs when a model is too simple or lacks the capacity to capture the underlying patterns in the data. It usually happens when the model is not complex enough to represent the true relationship between the input features and the target variable. The

most common reason for underfitting is using a linear model to represent a non-linear relationship in the data.

**What does it mean to overfit? When is it going to happen?**

ii. Overfitting occurs when a model becomes overly complex and fits the training data too closely. It happens when the model **captures noise or random fluctuations** in the training data, resulting in poor generalization to new, unseen data. Overfitting can occur when the model is too flexible or when there is insufficient regularization.

**In the sense of model fitting, explain the bias-variance trade-off.**

The bias-variance trade-off refers to the trade-off between a model's ability to capture the complexity of the data (low bias) and its tendency to be sensitive to noise or random fluctuations in the training data (high variance). Finding the right balance is crucial for building models that generalize well to new data.

**5. Is it possible to boost the efficiency of a learning model? If so, please clarify how.**

Yes, it is possible to improve the efficiency of a learning model. Some ways to boost model performance include:

- Increasing the amount and quality of training data.
- Feature engineering: Creating new informative features from existing ones or selecting the most relevant features.
- Tuning model hyperparameters: Optimizing the model's configuration to achieve better performance.
- Using ensemble methods: Combining multiple models to leverage their collective predictive power.
- Regularization: Applying techniques like L1 or L2 regularization to control model complexity and prevent overfitting.
- Cross-validation: Evaluating the model's performance on multiple subsets of the data to assess its generalization ability.

**6. How would you rate an unsupervised learning model's success? What are the most common success indicators for an unsupervised learning model?**

The success of an unsupervised learning model is typically measured by different indicators than those used for supervised learning. Common success indicators for unsupervised learning include:

- **Cluster quality metrics:** Assessing the cohesion (how close points within a cluster are) and separation (how distinct clusters are) of the discovered clusters.
- **Reconstruction error:** For dimensionality reduction techniques like autoencoders or PCA, measuring the quality of reconstructing the input data.
- **Visualization:** Examining the resulting representations or visualizations to assess if meaningful patterns or structures have been captured.

- **Domain-specific evaluation:** Depending on the application, domain-specific evaluation criteria can be used to determine the success of the unsupervised learning model.

**7. Is it possible to use a classification model for numerical data or a regression model for categorical data with a classification model? Explain your answer.**

Classification models are specifically designed to handle categorical data or predict categorical outcomes. They work by learning the relationships between input features and discrete class labels. Regression models, on the other hand, are used for predicting numerical or continuous values based on input features. While it is technically possible to misuse models (e.g., using a regression model for classification), it is generally recommended to use the appropriate model type based on the nature of the target variable.

**8. Describe the predictive modeling method for numerical values. What distinguishes it from categorical predictive modeling?**

Predictive modeling for numerical values, often referred to as regression, involves learning a function that maps input features to a continuous numerical output. The main difference from categorical predictive modeling is the type of output variable. Regression models aim to estimate or predict a specific numerical value, while categorical predictive models aim to assign instances to specific categories or classes.

**9. The following data were collected when using a classification model to predict the malignancy of a group of patients' tumors:**

- i). Accurate estimates – 15 cancerous, 75 benign
- ii). Wrong predictions – 3 cancerous, 7 benign

**Determine the model's error rate, Kappa value, sensitivity, precision, and F-measure.**

Based on the given data for tumor prediction:

- Error rate:  $(3 + 7) / (15 + 75) = 10 / 90 = 0.1111$  or 11.11%
- Kappa value: Calculating the Kappa value requires knowledge of the observed agreement and expected agreement. Without that information, it is not possible to calculate the Kappa value.
- Sensitivity:  $15 / (15 + 3) = 15 / 18 \approx 0.8333$  or 83.33%
- Precision:  $15 / (15 + 7) = 15 / 22 \approx 0.6818$  or 68.18%
- F-measure: The F-measure combines precision and recall and can be calculated using the formula:  $2 * (\text{precision} * \text{recall}) / (\text{precision} + \text{recall})$ .

**10. Make quick notes on:**

**1. The process of holding out:-**

The process of holding out refers to setting aside a portion of the available data, typically a validation or test set, to evaluate the model's performance on unseen data.

## 2. Cross-validation by tenfold

Cross-validation by tenfold is a technique where the data is divided into ten equal-sized subsets, called folds. The model is trained and evaluated ten times, with each fold used as a validation set once.

## 3. Adjusting the parameters

Adjusting the parameters involves tuning the hyperparameters of the model to optimize its performance. Hyperparameters are settings that are not learned from the data but are chosen by the practitioner.

## 11. Define the following terms:

### Purity vs. Silhouette width:

- Purity is a measure used in clustering to assess how well the instances within each cluster belong to the same class. It measures the homogeneity of clusters in terms of class labels.
- Silhouette width is a metric used to evaluate the quality of clustering results. It measures how well instances are clustered together and separated from other clusters based on the distances between them.

### Boosting vs. Bagging:

- Boosting is an ensemble method where **multiple weak learners (typically decision trees) are combined sequentially, with each learner** focusing on the instances that the previous learners struggled with. It aims to build a strong learner with improved predictive power.
- Bagging (Bootstrap Aggregating) is an ensemble method where **multiple independent learners are trained on different subsets of the training data**, often obtained through bootstrap sampling. The final prediction is usually made by combining the predictions of individual learners.

### The eager learner vs. the lazy learner:

- The eager learner (e.g., decision trees) constructs a generalized model during the training phase and uses it to make predictions at runtime. It eagerly learns from the training data and creates a representation of the learned knowledge.
- The lazy learner (e.g., k-nearest neighbors) postpones the generalization step during training and keeps the training instances stored. At runtime, it defers the learning process until a prediction is requested and then searches the stored instances for the most similar ones to make a prediction.