# Machine Learning #07

### ▼ 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?

In machine learning, *a target function is the function that a learning algorithm tries to approximate from the training data.* The target function maps inputs to outputs, and the goal of the learning algorithm is to produce a model that can make accurate predictions on new, unseen data.

*A real-life example of a target function could be predicting the price of a house based on its features such as location, number of bedrooms, square footage, etc.* The target function in this case would take these features as inputs and produce the predicted price as the output.

The fitness of a target function is typically assessed using a performance metric such as mean squared error or accuracy. The learning algorithm tries to minimize this metric by adjusting the model's parameters during training. Once the model is trained, it is evaluated on a separate test set to assess its performance on new, unseen data. The goal is to produce a model with a high level of fitness that can generalize well to new data.

### ▼ 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 and descriptive models are two types of models used in machine learning, each with its own purpose and characteristics.

*Predictive models are used to predict an outcome or behavior based on input data.* These models are designed to identify patterns in the input data and use those patterns to make predictions about future data points. *They are trained using historical data, where the outcome or behavior is known, and then used to make predictions on new, unseen data.*

*Examples of predictive models include:*

- Credit risk models that predict the likelihood of a borrower defaulting on a loan based on factors such as credit history, income, and debt-to-income ratio.

- Spam filters that predict the likelihood of an incoming email being spam based on the content of the email and other metadata.

- Stock price prediction models that predict the future price of a stock based on historical price data and other factors such as news and economic indicators.

*Descriptive models, on the other hand, are used to describe or summarize a dataset. These models are designed to identify patterns and relationships within the data and provide insights into the characteristics of the data.* They are used to answer questions such as "What happened?" and "What is happening?"

*Examples of descriptive models include:*

- Summary statistics such as mean, median, and standard deviation that describe the central tendency and variability of a dataset.

- Cluster analysis that groups data points into similar clusters based on their characteristics.

- Association rule mining that identifies relationships between different variables in a dataset.

In summary, *predictive models are used to make predictions based on input data, while descriptive models are used to describe or summarize a dataset. Predictive models are trained on historical data and used to make predictions on new, unseen data, while descriptive models are used to gain insights into the characteristics of a dataset.*

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

The efficiency of a classification model can be assessed using a variety of performance metrics. Some common performance metrics are:

1. *Confusion matrix:* A confusion matrix is a table that shows the true positive, false positive, true negative, and false negative values for a classification model. *It is a useful tool for evaluating the accuracy of a model's predictions.*

2. *Accuracy:* *Accuracy is the proportion of correct predictions made by the model.* It is calculated by dividing the number of correct predictions by the total number of predictions.

3. *Precision:* *Precision is the proportion of true positive predictions among all positive predictions.* It measures the model's ability to correctly identify positive cases.

4. *Recall:* *Recall is the proportion of true positive predictions among all actual positive cases.* It measures the model's ability to detect positive cases.

5. *F1 score:* *F1 score is the harmonic mean of precision and recall.* It is a single metric that balances both precision and recall.

6. *ROC curve:* The receiver operating characteristic (ROC) curve *is a plot of the true positive rate (sensitivity) against the false positive rate (1-specificity) at various thresholds. It is used to evaluate the performance of binary classification models.*

To evaluate the efficiency of a classification model, we typically divide the dataset into training and test sets. The model is trained on the training set, and its performance is evaluated on the test set using one or more of the above performance metrics.

For example, suppose we have a binary classification model that predicts whether a person will buy a product or not. We divide the dataset into a training set and a test set. We train the model on the training set and evaluate its performance on the test set using the confusion matrix, accuracy, precision, recall, F1 score, and ROC curve.

The confusion matrix shows the number of true positive, false positive, true negative, and false negative values. The accuracy measures the proportion of correct predictions made by the model. The precision measures the proportion of true positive predictions among all positive predictions, and recall measures the proportion of true positive predictions among all actual positive cases. The F1 score is the harmonic mean of precision and recall. Finally, the ROC curve plots the true positive rate against the false positive rate at various thresholds.

*By evaluating the model's performance using multiple performance metrics, we can gain a comprehensive understanding of the model's strengths and weaknesses and identify areas for improvement.*

## ▼ 4. Describe:

i. In the sense of machine learning models, what is underfitting? What is the most common reason for underfitting?

ii. What does it mean to overfit? When is it going to happen?

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

i. *Underfitting* in machine learning models refers to a scenario **where the model's complexity is low, and it <u>fails to capture the underlying patterns</u> in the data.** It can occur due to various reasons, such as using a model that is too simple or having insufficient data to train the model. The most common reason for underfitting is using a model that is too simple for the given dataset.

ii. *Overfitting* in machine learning models refers to a scenario **where the model is too complex and fits the training data too closely, resulting in <u>poor generalization to new, unseen data</u>.** It typically happens when the model has too many parameters relative to the amount of data available for training. As a result, the model may start to memorize the training data instead of learning the underlying patterns.

iii. *The bias-variance trade-off* is a fundamental concept in model fitting that **refers to the relationship between a model's ability to fit the training data (bias) and its ability to generalize to new, unseen data (variance).** *A high bias model is one that is too simple and may underfit the data, whereas a high variance model is one that is too complex and may overfit the data.* The goal is to find a model that strikes a balance between bias and variance and generalizes well to new data.

*In general, increasing a model's complexity can reduce its bias but increase its variance, while decreasing its complexity can increase its bias but decrease its variance.* Therefore, finding the optimal balance between bias and variance is crucial to building an accurate and generalizable machine learning model. This is typically achieved through techniques such as regularization and cross-validation.

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

<u>*Yes*</u>*, it is possible to boost the efficiency of a learning model*. Some of the ways to achieve this are:

1. *Feature Engineering*: Feature engineering **<u>involves selecting or transforming the input variables</u>** to enhance the predictive power of the model. It **<u>helps to extract the most relevant features from the data and reduces the dimensionality</u>** of the input space. **<u>This can result in faster training and better performance of the model.</u>**

2. *Hyperparameter Tuning:* **Hyperparameters are model parameters that are <u>set before training</u> and can significantly impact the model's performance.** Tuning these hyperparameters **<u>can improve the model's efficiency and accuracy.</u>** Techniques such as grid search, random search, and Bayesian optimization can be used to find the optimal values for these hyperparameters.

3. *Regularization:* **Regularization is a technique that helps <u>prevent overfitting</u> by <u>adding a penalty</u> term to the loss function.** This can reduce the model's complexity and improve its generalization ability.

4. *Ensemble Learning:* **Ensemble learning involves combining multiple models to improve the overall performance of the model.** Techniques such as bagging, boosting, and stacking can be

used to achieve this.

5. *Transfer Learning:* **Transfer learning is a technique that involves <u>leveraging pre-trained models</u> for a similar task to improve the efficiency of the model.** This can help reduce the training time and improve the performance of the model.

Overall, boosting the efficiency of a learning model requires careful analysis and optimization of various parameters and techniques.

## ▼ 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 can be evaluated using various metrics, depending on the specific task and the type of model being used. Here are some common success indicators for unsupervised learning models:

1. *Clustering Quality:* In unsupervised learning, clustering is a common task that involves grouping similar data points together. **The quality of the clustering can be evaluated using metrics such as silhouette score, Dunn index.** These metrics measure how well the data points are separated into distinct clusters.

2. *Reconstruction Error:* For models **that involve data compression or dimensionality <u>reduction</u>**, such as PCA or autoencoders, the reconstruction error is an important metric to evaluate the model's success. **The reconstruction error measures how well the model can reconstruct the original data from the compressed or reduced version.**

3. *Anomaly Detection:* **In unsupervised anomaly detection, the goal is to identify data points that deviate significantly from the norm**. The success of the model can be evaluated using metrics such as precision, recall, F1-score, or area under the receiver operating characteristic curve (AUROC).

4. *Visualization:* Unsupervised learning models can be used for visualization, where high-dimensional data is projected onto a lower-dimensional space for better understanding. **The success of the model can be evaluated by how well the data is separated or clustered in the visualization.**

Overall, the success of an unsupervised learning model depends on the specific task and the evaluation metrics chosen. It is important to carefully choose appropriate metrics that capture the model's performance and assess its success.

## ▼ 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.

*<u>No</u>, it is not appropriate to use a classification model for numerical data or a regression model for categorical data because these models are designed to handle different types of data.*

*Classification models are used to predict categorical outcomes, where the <u>output variable</u> can take on a <u>limited number of discrete values</u>.* The goal of classification is to assign a class label to each input data point. On the other hand, *regression models are used to predict continuous outcomes, where the output variable can take on any numeric value within a given range*. The goal of regression is to find a function that predicts the value of the output variable based on the input data.

Trying to use a classification model for numerical data or a regression model for categorical data will lead to incorrect predictions and poor model performance. For example, if we use a classification

model to predict a numeric value, the model will only be able to output a limited number of discrete values, which will not accurately represent the underlying numeric value. Similarly, if we use a regression model to predict a categorical variable, the model will output a continuous value, which cannot be mapped to the limited set of categorical values.

Therefore, it is important to choose the appropriate model type based on the nature of the data and the problem being solved.

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

*Predictive modeling for numerical values is also known as regression modeling. It is a method used in machine learning to build a model that can predict a continuous numeric value for a given set of input features.* **The goal of regression modeling is to find a mathematical function that best describes the relationship between the input variables and the output variable.**

Regression models can be used to predict a wide range of numerical outcomes, such as stock prices, housing prices, or the number of sales. There are various types of regression models, such as linear regression, polynomial regression, and logistic regression.

On the other hand, *categorical predictive modeling is used to predict a categorical or discrete outcome variable, where the output variable can take on a limited number of discrete values.* **The goal of classification modeling is to assign a class label to each input data point.** Classification models can be used to predict outcomes such as spam or not spam, fraudulent or not fraudulent, or customer churn or no churn.

*The key difference between numerical and categorical predictive modeling is the type of output variable that is being predicted.* Numerical predictive modeling is used for continuous numeric outcomes, while categorical predictive modeling is used for discrete categorical outcomes. As a result, different modeling techniques and evaluation metrics are used for each type of modeling.

## ▼ 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.

Using the given information, we can calculate the various performance metrics for the classification model as follows:

- True Positives (TP) = 15 (cancerous patients correctly classified)

- False Positives (FP) = 7 (benign patients incorrectly classified as cancerous)

- False Negatives (FN) = 3 (cancerous patients incorrectly classified as benign)

- True Negatives (TN) = 75 (benign patients correctly classified)

Error rate = (FP + FN) / (TP + TN + FP + FN) = (7 + 3) / (15 + 75 + 7 + 3) = 10 / 100 = 0.1 or 10%

Kappa value = (accuracy - chance agreement) / (1 - chance agreement) = (0.9 - 0.65) / (1 - 0.65) = 0.5

Sensitivity = TP / (TP + FN) = 15 / (15 + 3) = 0.833 or 83.3%

Precision = TP / (TP + FP) = 15 / (15 + 7) = 0.682 or 68.2%

F-measure = 2 * (precision * sensitivity) / (precision + sensitivity) = 2 * (0.682 * 0.833) / (0.682 + 0.833) = 0.750 or 75.0%

Therefore, the error rate of the model is 10%, the Kappa value is 0.5, the sensitivity is 83.3%, the precision is 68.2%, and the F-measure is 75.0%.

## ▼ 10. Make quick notes on:

1. The process of holding out

2. Cross-validation by tenfold

3. Adjusting the parameters

1. *The process of holding out: A technique used in machine learning to <u>evaluate the performance</u> of a model <u>on an independent dataset.</u>* A portion of the dataset is held out, or reserved, and not used during the training process. The model is then evaluated on this held-out dataset to test its generalization ability.

2. *Cross-validation by tenfold: A technique used in machine learning to evaluate the performance of a model by splitting the dataset into ten equal parts or folds.* The model is trained on nine of the folds and tested on the remaining fold. This process is repeated ten times, with each fold being used once for testing.

3. *Adjusting the parameters: A process in machine learning that involves tuning the parameters of a model to achieve better performance.* This is typically **done by changing the values of hyperparameters,** which are set before training the model. The goal is to find the optimal combination of hyperparameters that results in the best performance of the model.

## ▼ 11. Define the following terms:

1. Purity vs. Silhouette width

2. Boosting vs. Bagging

3. The eager learner vs. the lazy learner

1. *Purity vs. Silhouette width*:

   - *Purity:* A measure used in cluster analysis to evaluate the quality of clustering results. *It measures how well each cluster contains only a single class of data points.*

   - *Silhouette width*: A measure used in cluster analysis to evaluate the quality of clustering results. *It measures how similar each data point is to its own cluster compared to other clusters.*

2. *Boosting vs. Bagging:*

   - *Boosting: An ensemble learning technique that <u>combines multiple weak models to create a strong model.</u>* The models are trained sequentially, with each subsequent model trying to correct the errors of the previous models.

   - *Bagging: An ensemble learning technique that <u>combines multiple models by training them on different subsets of the training data.</u>* The models are trained independently, and their predictions are combined to make the final prediction.

3. *The eager learner vs. the lazy learner:*

   - *Eager learner: A machine learning algorithm that builds a model <u>before receiving any data,</u> and then uses the model to make predictions when new data is introduced.* Eager learners are also called "eager classifiers" because they are eager to classify new data.

- ***Lazy learner:*** **A machine learning algorithm that postpones building a model until it is given new data to classify.** Lazy learners are also called "lazy classifiers" because they are lazy to build a model. They make predictions by comparing the new data to the stored training data.