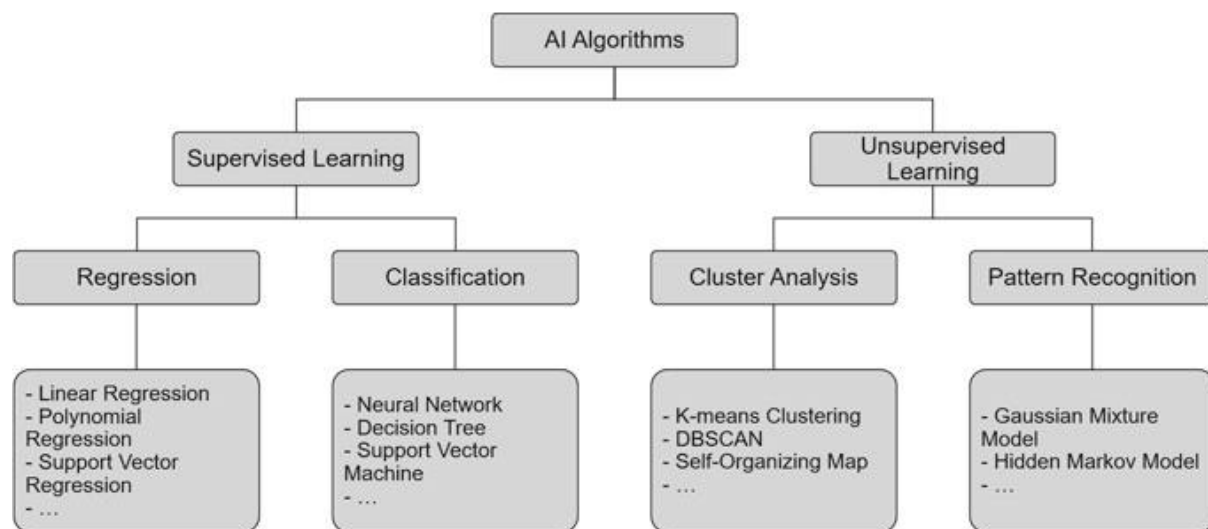# Industrial AI: Categories of Algorithms

Functionally, AI algorithms can be divided into the categories of supervised learning and unsupervised learning according to whether or not the labels of data are known during training. As shown in Fig, supervised learning refers to how the data input during model training includes input objects (usually a set of feature vectors) and their expected values (discrete or continuous values). According to the training, supervised learning aims at generating the relationship between input feature vectors and expected values, generating a mathematical model with an inferential prediction



function for mapping new input objects. During training, supervised learning can be classified into regression or classification, according to the continuity of expected values.

In regression algorithms, the training labels are continuous, and the trained model can infer the corresponding expected value according to the new input object. In the classification algorithm, the training labels are discrete values, and the trained model can classify the unlabeled new input objects into corresponding categories.

Unlike supervised learning, training data in unsupervised learning does not contain labels. The process of learning is not to find the relationship between the input object and the expected labels, but to recognize the patterns amongst the input objects. Typical unsupervised learning processes include clustering algorithms and statistical estimation. Clustering algorithms aim at grouping similar input objects together to form distinct categories, while statistical estimation algorithms use the principle of probability and statistics to describe the input object as a distribution function with specific parameters or to estimate the correlation of a time series between input data.In the following section, we will introduce several typical learning models for each AI algorithm and their applications in industrial data analysis.

### *Regression Algorithms*

In regression algorithms, the input feature vector $\_x$ and the label value $y$ in the training dataset are used to establish an estimation function so that $y = f(x)$ or $y \approx f(x)$. When such an estimation function or model is established, we can estimate the expected value $y$ of the new input feature vector $x$. In industrial scenarios, such algorithms are often used for virtual metrology or system health assessment. Virtual metrology means that the quality of the product or stability of the production process on the line can be estimated by the data collected in the production process, without additional measurement or detection processes. The health assessment of a system usually uses the input and output of the system to establish amodel.The system output obtained by the comparative measurement and the system output estimated by the model are used to achieve the monitoring. The regression algorithms commonly used in the industrial domain are linear regression, polynomial regression, and support vector regression.

**Classification Algorithms**

Similar to regression algorithms, the classification algorithm also wants to obtain an estimation function f (·), but the result of the estimation is not a continuous value; it is a discrete value. Possible discrete values form a set in which each element is a possible category. The task of the classification algorithm is to construct such a classifier, f (·), so that there is only one class corresponding to the input feature vector.

In industrial data analysis, classification algorithms are often used to diagnose faults or trace their cause. Classification models are established by using historical data and the corresponding labels. In the process of analyzing the newly collected data, the classification model can estimate the possible failure modes or the important factors that cause certain types of faults. The classification algorithms commonly used in the industrial domain mainly include support vector machines, neural networks, and decision trees.

**Clustering Algorithms**

Clustering algorithms cluster the input feature vectors unsupervised on the basis of similarity, which makes the similarity within clusters larger than that between clusters. Because input data has no label information in the clustering process, the clustering only estimates the dataset and recognizes the pattern. Clustering results often need expert experience to explain the characteristics of each cluster. In industrial scenarios, clustering algorithms are mainly used to identify different working conditions and to evaluate the health of the system.

When the working conditions of the system are complex, it is necessary to identify these different conditions before analyzing the data in order to establish corresponding analysis methods for different working conditions. In the absence of labelled historical data, we can first cluster the data with its health status to

identify patterns. The health status of the system can be obtained by comparing the newly collected data with the identified health patterns. Common clustering algorithms include K-means cluster, DBSCAN, and self-organizing maps.

**Statistical Estimation Algorithms**

Statistical estimation is another standard unsupervised pattern recognition algorithm. It uses the principles of statistics and probability theory to identify the potential statistical distribution form of the input dataset for estimation. Datasets can be represented as a combination of one or more distribution functions—the relationship of a time series between each input feature vector of data can also be represented by probability and state transition functions.

Statistical estimation can be used to represent the current state of the system, or the entire decay process of the system, and can characterize the distribution of data under different fault modes and degrees. Based on estimated data distribution patterns, users can deepen their understanding and quantify the risks of equipment operation. Common statistical estimation algorithms are the hidden Markov model and the Gaussian mixed model.