

Chapter 6

Artificial Intelligence for Dynamic Spectrum Management



Abstract In the past decade, a significant advancement has been made in artificial intelligence (AI) research from both theoretical and application perspectives. Researchers have also applied AI techniques, particularly machine learning (ML) algorithms, to DSM, the results of which have shown superior performance as compared to traditional ones. In this chapter, we first provide a brief review on ML techniques. Then we introduce recent applications of ML algorithms to enablers of DSM, which include spectrum sensing, signal classification and dynamic spectrum access.

6.1 Introduction

Artificial intelligence (AI), also known as machine intelligence, has been seen as the key power to drive the development of future information industry [1]. The term AI was coined by John McCarthy in a workshop at Dartmouth College in 1956, and he defined AI as “the science and engineering of making machines, especially intelligent computers” [2]. Generally, AI is defined as the study of the intelligent agent, which is able to judge and execute actions by observing the surrounding environment so as to complete certain tasks. The intelligent agent can be a system or a computer program. With the significant advancement in the computational capability of computer hardware, various theories, especially machine learning techniques, and applications of AI have been developed in the past two decades.

With the surging demand for wireless services and the increasing connections of wireless devices, the network environments are becoming more and more complex and dynamic, which imposes stringent requirements on DSM. In the age of 5G, AI has been seen as an effective tool to support DSM in order to tackle the transmission challenges, such as high rate, massive connections and low latency [3, 4]. By adopting ML techniques, the traditional model-based DSM schemes would be transformed to the data-driven DSM schemes, in which the controller in the network can adjust itself adaptively and intelligently to improve the efficiency and robustness of DSM. AI-based DSM schemes have thus attracted more and more attention in recent years, and have shown great potentials in practical scenarios.

The applications of AI techniques would bring significant benefits to DSM. Firstly, the AI-based DSM schemes normally do not need environmental information as the prior knowledge, and they can extract useful features from the surroundings automatically. Secondly, AI-based DSM schemes can be re-trained periodically and thus they are more robust to the changing environment. Additionally, by applying AI techniques, DSM can be done in a decentralized and distributed manner, leading to a significant reduction of signal overheads, especially for large-scale systems. Finally, once trained, AI-based DSM schemes are low in complexity for processing newly arrived data and thus they are more suitable for practical implementation.

While it is believed that machine learning techniques are effective methods for developing and optimizing the next generation networks [5], there also exist some challenges in applying AI techniques in DSM. For example, different from images, the received signal and its higher order statistics in wireless networks are normally complex numbers, which are hard to process directly by neural networks. Additionally, in a typical wireless communication system, accurate network data such as channel information, is hard to obtain in practice. Hence, there are many remaining challenges and problems to be addressed for achieving wireless intelligence. In this chapter, we first provide a brief review of machine learning techniques, then introduce some applications of these algorithms to DSM, including spectrum sensing, signal classification and dynamic spectrum access.

6.2 Overview of Machine Learning Techniques

As the core technique of AI, machine learning (ML) is a multidisciplinary subject involving multiple disciplines such as probability theory, statistics, information theory, computational theory, optimization theory, and computer science. T. Mitchell provided a brief definition of machine learning in 1997 as follows: “machine learning is the study of computer algorithms that improve automatically through experience” [6]. Hence, the main objective of ML is to make agents simulate or implement human learning behaviors. For example, with the help of ML algorithms, a machine agent is able to learn from training data to achieve different tasks such as image recognition.

Based on the type of training data used, ML can be divided into two branches, namely, supervised learning and unsupervised learning. The former requires labeled training data, while the latter only uses unlabeled training data.

In supervised learning, the objective for an agent is to learn a parameterized function from the given labeled training dataset and then based on the function learnt to predict the result directly while new data arrives. The common tasks in supervised learning are regression and classification. Specifically, regression is to determine the quantitative relationship between certain variables based on a set of training data, and classification is to find a function to determine the category to which the input data belongs.

In unsupervised learning, since the training data is unlabeled, the agent needs to adopt clustering methods to obtain the relationship. A clustering method aims to

divide the training data into several classes based on the similarity of the data. The objective of clustering is to minimize intra-class distance while maximizing inter-class distance. Compared to supervised learning, unsupervised learning is more like self-study.

Labeled data can also be generated through online learning such as reinforcement learning (RL). In particular, RL produces labeled experiences to train itself from continuous interactions with the environment, it is developed to solve a *Markov decision process* (MDP) $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}\}$, where \mathcal{S} is the state space, \mathcal{A} is the action space, \mathcal{P} is the transition probability space and \mathcal{R} is the reward function [7].

ML techniques can also be grouped into two categories, namely, statistical machine learning (SML) and deep learning (DL). Using statistics and optimization theory, SML constructs proper probabilistic and statistical models with training data. DL, on the other hand, makes use of artificial neural network (ANN), also known as deep neural network (DNN), to perform supervised learning tasks. In recent years, neural network techniques have also been applied to RL, leading to the birth of deep reinforcement learning (DRL). In the following, we will provide a brief introduction to SML, DL and DRL.

6.2.1 Statistical Machine Learning

The objective of SML is to construct a probabilistic and statistical model using the training data, then, based on the constructed model, to make inferences with new data [8]. SML can be applied in both supervised learning and unsupervised learning. The commonly used supervised learning methods with SML are support vector machine (SVM) and K-nearest neighbor (KNN), and the commonly used unsupervised learning methods with SML are K-means and Gaussian mixture model (GMM).

1. **K-Nearest Neighbor:** K-Nearest Neighbor (KNN) algorithm is a basic supervised learning algorithm for classification. Let $T = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$ denote a given training dataset, where \mathbf{x}_i is the i -th data set and y_i is the corresponding label. Assume that all data sets come from J classes. For a newly arrived data set \mathbf{x} , its label, i.e., class, is determined by its K nearest labeled neighbors based on the adopted classification decision rules. Hence, the basic elements in KNN are the number of neighbors K , the distance measure and the classification decision rule.

Specifically, the classification process consists of two steps: the first step is to search K labeled data sets which are closest to the newly arrived data set \mathbf{x} according to the given distance measure. Denote the region covering these K data sets as $N_K(\mathbf{x})$. The second step is to determine its label y by using the chosen classification decision rule based on $N_K(\mathbf{x})$. The commonly used classification decision rule is the majority voting rule, which is given as

$$y = \arg \max_{c_j, j=1, \dots, J} \sum_{\mathbf{x}_i \in N_K(\mathbf{x})} I(y_i = c_j) \quad (6.1)$$

where $I(\cdot)$ is the indicator function that indicates if a label belongs to class c_j . For example, $I(y_i = c_j) = 1$ if $y_i = c_j$, and $I(y_i = c_j) = 0$, otherwise.

2. **Support Vector Machine:** Support vector machine (SVM) algorithm is a typical binary classification algorithm. The basic idea of the SVM algorithm is to find a decision hyperplane to maximize the margin between different classes. Specifically, for a given training data set $T = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$, where $y_i \in \{-1, 1\}$, the objective of the SVM algorithm is to find the hyperplane denoted by $\mathbf{w} \cdot \mathbf{x} + b = 0$ to make the data sets linearly separable, where \mathbf{w} and b are the normal vector and the intercept of the plane, respectively. If the decision hyperplane is obtained, the corresponding classification decision function is given as

$$f(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b) \quad (6.2)$$

The hyperplane can be learnt by solving the following convex quadratic programming problem

$$\min \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i \quad (6.3)$$

$$\text{s.t. } y_i (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i, \quad i = 1, 2, \dots, N \quad (6.4)$$

$$\mathbf{x}_i \geq 0, \quad i = 1, 2, \dots, N \quad (6.5)$$

where C is a punishment parameter and ξ_i is the soft constant for i -th data set. Generally, SVM is used to solve a linear classification problem, but it can also be used as a nonlinear classifier by introducing different kernel functions such as Gaussian kernel function and radial basis function.

3. **K-means:** K-means algorithm is a clustering algorithm, in which the unlabeled data sets are processed iteratively to form K clusters. Specifically, at the beginning, K data sets are chosen to form the initial centroids of the K clusters. Then, the K-means algorithm alternates the following two steps. The first step is to assign each of the remaining data sets to its nearest cluster. This is determined by evaluating the Euclidian distance between the data set and the centroid of each cluster and choosing the cluster with the smallest distance. The second step is to update the centroid of each cluster, denoted as c_k , based on the newly labeled data sets. Mathematically, this can be expressed as

$$c_k = \frac{1}{|N_k|} \sum_{\mathbf{x} \in N_k} \mathbf{x}, \quad k = 1, 2, \dots, K \quad (6.6)$$

where N_k denote the set of examples assigned to cluster k . These two steps will repeat until a termination condition is met.

Although K-means algorithm can be implemented with low complexity, its performance is influenced significantly by the initialization parameters such as the number of clusters and the cluster centroids.

4. **Gaussian Mixture Model:** Gaussian mixture model (GMM) is a widely used model for unsupervised learning. The probability density function (PDF) of the data sets can be expressed as

$$p(\mathbf{x}_i) = \sum_{k=1}^K \pi_k p(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (6.7)$$

where K is the number of Gaussian components, π_k is the mixing coefficient that satisfies $\sum_{k=1}^K \pi_k = 1$, and $p(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ denotes the PDF of the k th Gaussian component with mean $\boldsymbol{\mu}_k$ and covariance $\boldsymbol{\Sigma}_k$, which can be expressed as

$$p(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \frac{1}{\pi |\boldsymbol{\Sigma}_k|} \exp\left(-(\mathbf{x}_i - \boldsymbol{\mu}_k)^H \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)\right) \quad (6.8)$$

The unknown parameters of the GMM can be denoted as $\Theta = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$. The objective of the GMM algorithm is to find the optimal parameters $\Theta = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$ to maximize the following log-likelihood function

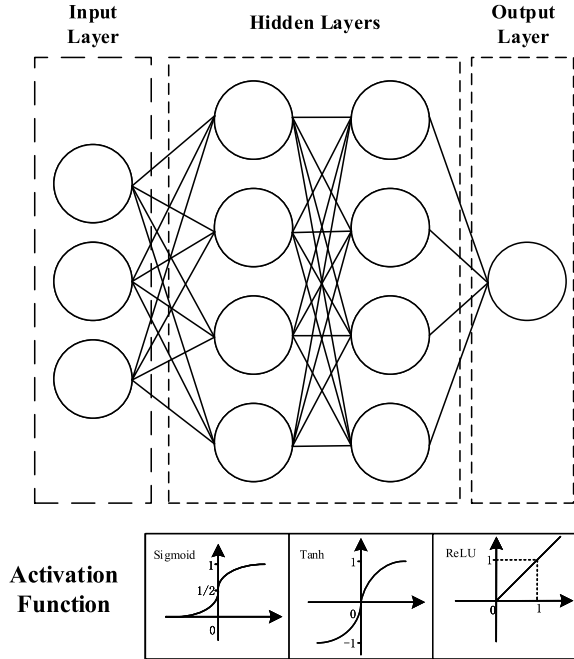
$$\mathcal{L}(\Theta) = \sum_{i=1}^N \ln \left(\sum_{k=1}^K \pi_k \cdot p(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right) \quad (6.9)$$

Since there is no closed-form solution for the above problem, the expectation maximization (EM) algorithm is usually adopted to solve for the optimal parameters $\Theta = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$ in an iterative manner with properly chosen initial values for them. The EM algorithm is generally composed of two steps in each iteration, namely, the expectation step and the maximization step. Denote γ_{ik} as a latent variable which represents the probability that example \mathbf{x}_i belongs to the k -th cluster. In the expectation step, the latent variable γ_{ik} is updated as: $\gamma_{ik} = \frac{\pi_k p(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{k=1}^K \pi_k p(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}$, for $i = 1, \dots, N$ and $k = 1, \dots, K$. In the max-

imization step, the parameter Θ is updated as: $\pi_k = \frac{1}{N} \sum_{i=1}^N \gamma_{ik}$, $\boldsymbol{\mu}_k = \frac{\sum_{i=1}^N \gamma_{ik} \mathbf{x}_i}{\sum_{i=1}^N \gamma_{ik}}$, and

$$\boldsymbol{\Sigma}_k = \frac{\sum_{i=1}^N \gamma_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^H}{\sum_{i=1}^N \gamma_{ik}}, \text{ for } k = 1, \dots, K.$$

Fig. 6.1 The basic model of a fully-connected ANN, which is composed of input layer, output layer and hidden layers. The commonly used activation functions are sigmoid function, Tanh function and ReLU function



6.2.2 Deep Learning

Deep learning (DL) has significantly advanced the development of computer vision (CV) and natural language processing (NLP) recently. As the core technique of DL, ANN has been used to approximate the relationship between an input and an output. Generally, a typical ANN is composed of three parts, namely, input layer, output layer and hidden layers as shown in Fig. 6.1. In each layer, many cells with different activation functions are placed, and the cells in adjacent layers are connected with each other in a pre-designed manner. With the development of ANNs, there are different network structures used for different types of data. For example, a convolutional neural network (CNN), which consists of convolutional layers, pooling layers and fully connected layers, is suitable for images; while a recurrent neural network (RNN), which contains many recurrent cells in the hidden layers, is suitable for time series data. Furthermore, in order to improve the generalization and convergence performance of the DL, dropout and other techniques are introduced in the design of neural networks [9].

1. **Convolutional Neural Network:** Convolutional neural network (CNN) is a special network for processing images, in which the cells adopt convolution operations. A typical CNN is composed of multiple convolutional layers, pooling layers and fully-connected layers [10].

- a. **Convolutional Layer:** Different from the fully-connected layers, a convolutional layer contains a set of multiple feature maps, which are obtained by using different convolution kernels to operate on the input image. In particular, one feature map is calculated by one convolution kernel operating on the input image, which means all the elements in the same feature map share the same weight and bias with each other. Besides, the size of the convolution kernel is smaller than that of the input, and the CNN has the unique characteristic referred to as *sparse interactions*.
- b. **Pooling Layer:** A pooling layer is usually placed after a convolutional layer in a CNN to capture invariant features. Specifically, the pooling operation is to replace the output of one position in the input image with a summary statistic of the neighborhood. A commonly used pooling method is the max-pooling function, which gives the maximum output of the rectangular region. In fact, the pooling operation can be seen as an action to add a strong prior knowledge.

In order to utilize the extracted feature maps, fully-connected layers are normally used as the last several layers of a CNN. With the help of the special structure, CNNs can process data with clear mesh topology effectively.

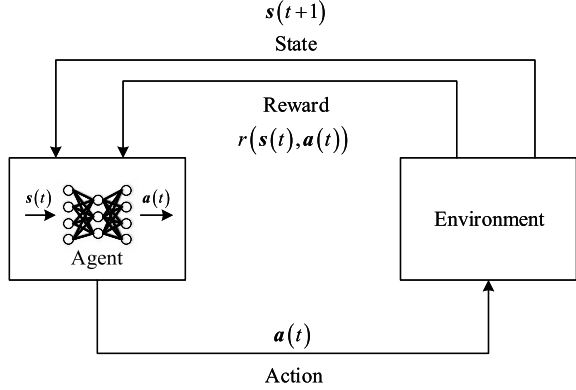
2. **Recurrent Neural Network:** Recurrent neural network (RNN) is a powerful tool for time series data, which have shown superior performance on speech recognition [11]. Different from traditional neural networks, there are many connected cells in each layer in an RNN. All cells in the same layer have the same structure and each of them passes its information to its successor. The output of an RNN is determined by not only its current input but also the memory recorded in the past time steps. However, conventional RNNs cannot learn long-term dependent information and suffer from the gradient vanishing problem easily. Then long short-term memory (LSTM) network, as a kind of gated RNN network, is proposed to mitigate this problem. Specifically, in each cell of an LSTM network, there are three gates, namely, the input gate, the forget gate and the output gate, which are given as follows

$$\begin{aligned}
 \mathbf{i}_t &= \sigma(\mathbf{W}_i \mathbf{h}_{t-1} + \mathbf{U}_i x_t + b_i) \\
 \mathbf{f}_t &= \sigma(\mathbf{W}_f \mathbf{h}_{t-1} + \mathbf{U}_f x_t + b_f) \\
 \mathbf{o}_t &= \sigma(\mathbf{W}_o \mathbf{h}_{t-1} + \mathbf{U}_o x_t + b_o)
 \end{aligned} \tag{6.10}$$

where \mathbf{i}_t , \mathbf{f}_t , and \mathbf{o}_t are the input gate, the forget gate and the output gate, respectively; \mathbf{W}_i , \mathbf{U}_i , b_i , \mathbf{W}_f , \mathbf{U}_f , b_f , \mathbf{W}_o , \mathbf{U}_o , b_o are the weight matrices and biases of the corresponding gate, respectively; and $\sigma(\cdot)$ is the sigmoid function. Additionally, each cell has a self-loop and its cell state is jointly controlled by the forget gate and the input gate. Specifically, the forget gate determines what information to remove and the input gate determines what information to add to the next cell state. Mathematically, the cell state can be expressed as

$$\mathbf{c}_t = \mathbf{f}_t \cdot \mathbf{c}_{t-1} + \mathbf{i}_t \cdot \tanh(\mathbf{W}_c \mathbf{h}_{t-1} + \mathbf{U}_c x_t + b_c) \tag{6.11}$$

Fig. 6.2 The standard DRL framework, in which an agent selects actions by inputting its state to the neural network and executes the action to interact with the environment continuously



where W_c , U_c and b_c are the weight matrices and the bias of the cell memory, respectively. The gated structure allows the LSTM network to learn the long-term dependent information while avoiding vanishing gradients.

6.2.3 Deep Reinforcement Learning

As the combination of DL and RL, deep reinforcement learning (DRL) has shown superior performance in sequential decision-making tasks. In the DRL framework, as shown in Fig. 6.2, the agent inputs its observation (state) $s(t) \in \mathcal{S}$ into the neural network and outputs an action $a(t) \in \mathcal{A}$. Then it obtains a reward $\mathcal{R}(s(t), a(t))$ which is used to evaluate the profit of the selected action by executing it. After a period of learning, the agent can learn the optimal strategy, which maps a state to an action, to maximize its long-term accumulative reward from continuous interactions with the environment. Similar to RL, the basic elements of DRL are also state space \mathcal{S} , action space \mathcal{A} and the reward function \mathcal{R} . Different from the traditional RL, which uses a table to indicate the relationship between the state space and the action space, DRL uses a neural network as the function approximator, and therefore it works more effectively for problems with high dimensional state and action spaces. In DRL, the commonly used methods are deep Q-network (DQN), double deep Q-network (DDQN), asynchronous advantage actor-critic (A3C) and deep deterministic policy gradient (DDPG).

1. **Deep Q-network:** Different from the tabular method in traditional RL, a neural network called deep Q-network (DQN) is adopted to approximate the relationship between state space and action space [12]. Since the DQN is optimized by minimizing the temporal difference error, the loss function of DQN is given as

$$L(\theta) = \mathbb{E} \left[\left(y^{DQN} - Q(s, a; \theta) \right)^2 \right] \quad (6.12)$$

where $\mathbb{E}[\cdot]$ indicates the expectation operation, $Q(s, a; \theta)$ is the Q-function with the parameter θ , and the target value y^{DQN} is given as

$$y^{DQN} = \mathcal{R}(s, a) + \gamma Q(s', a'; \theta) \quad (6.13)$$

To improve the performance of the basic DQN, two other techniques, i.e., experience replay and quasi-static target network, are introduced in the design of DQN technique.

- **Experience Replay:** The agent needs to construct a fixed-length memory \mathcal{M} , which is based on the *first-in-first-out* (FIFO) rule. In each training step t , the agent needs to store newly obtained experience into the memory \mathcal{M} , and then a mini-batch d_t of experiences is sampled randomly from \mathcal{M} for training.
- **Quasi-static Target Network:** The agent constructs two DQNs of the same structure, i.e., the target DQN $Q(s, a; \theta')$ and the trained DQN $Q(s, a; \theta)$, where θ' and θ are their respective parameters. In every K steps, the trained network share its parameter with the target network.

Additionally, in order to balance the relationship between exploration and exploitation, the ϵ -greedy algorithm is usually adopted in a DRL. Specifically, an agent selects the action corresponding to the maximum Q-value of the trained network with a probability $1 - \epsilon$, and selects an action randomly otherwise. After the algorithm converges, the agent just selects the action with the maximum Q-value, and the target network is closed.

2. **Double Deep Q-network:** Since the target value is from the same DQN, the Q-function may be overestimated and trapped in a local optimum, leading to the performance degradation. To improve the performance of DQN, double deep Q-network (DDQN) can be adopted to provide more accurate estimation of the Q-function [13]. In the DDQN, the target value y^{DDQN} can be expressed as follows

$$y^{DDQN} = \mathcal{R}(s, a) + \gamma Q\left(s', \arg \max_{a' \in \mathcal{A}} Q(s', a'; \theta); \theta'\right) \quad (6.14)$$

After years of development, ML has become the most concerned discipline in the information age and shown strong effectiveness in applications. As the main force of AI technique, more and more ML algorithms are applied in various fields in order to achieve industrial intelligence.

6.3 Machine Learning for Spectrum Sensing

Spectrum sensing is an important task to realize DSM in wireless communication systems, and is usually used to assist users to find out the channel status. In order to increase the accuracy of spectrum sensing, many spectrum sensing algorithms have been developed in the past years, such as estimator-correlator (EC) detector, the semi-

blind energy detector and the blindly combined energy detection (BCED). Although the EC detector can achieve the optimal performance, it needs the knowledge of PU signals and noise level. The semi-blind energy detector is more practical, and it only requires the knowledge of the noise power. However, the performance of the semi-blind energy detector depends heavily on the accurate knowledge of noise power, which is usually uncertain. The BCED does not need any prior knowledge about the PU signals or noise, but the performance is worse than the performance of the semi-blind energy detector. It is noticed that most existing algorithms are model-driven, and need the prior knowledge of noise or PU signals to achieve good performance. However, this feature makes them unsuitable for practical environment, and the lack of prior knowledge would result in performance degradation.

To solve the above issues, machine learning techniques have been adopted to develop cooperative spectrum sensing (CSS) framework [14]. Specifically, the work considers a CR network, in which multiple SUs share a frequency channel with multiple PUs. The channel is considered to be unavailable for SUs to access if at least one PU is active and it is available if there is no active PU. For cooperative sensing, each SU estimates the energy level of the received signals and reports it to another SU who acts as a fusion center. After the reports of the energy level from all SUs are collected, the fusion center makes the final classification of the channel availability.

Using the machine learning technique such as K-means algorithm, GMM clustering, SVM algorithm and KNN algorithm, the fusion center can construct a classifier to detect the channel availability. With unsupervised machine learning such as K-means and GMM clustering, the detection of the channel availability relies on the cluster that the sensing reports from all the SUs are mapped to. On the other hand, with supervised machine learning such as SVM algorithm and KNN algorithm, the classifier is first trained using the labeled sensing reports from all SUs. After the classifier is trained, it can be directly used to derive the channel availability. Compared with traditional CSS techniques, the proposed machine learning framework can bring the following two advantages: (1) it is robust to the changes in the radio environment; (2) it can achieve a better performance in terms of classification accuracy.

6.4 Machine Learning for Signal Classification

Signal classification, usually performed before signal detection, is a fundamental task in cognitive radio networks. Consider the modulation classification as an example. Traditionally, there are two kinds of modulation classification approaches, namely, the likelihood-based (LB) approach and the feature-based (FB) approach. The LB approach is based on computing the likelihood function of received signals under different modulation schemes hypotheses, and the modulation scheme with the maximum likelihood value is validated. With perfect knowledge of channel and noise parameters, the LB approach can achieve the optimal performance in a Bayesian sense. However, the estimation of these parameters imposes high compu-

tation complexity. In the FB approach, useful features such as higher-order statistics are extracted for decision-making. In general, the FB approach has lower computational complexity but it can only achieve sub-optimal performance. Therefore, in order to achieve near optimal performance with low computational complexity, ML techniques have been introduced in solving the modulation classification problem, and have shown superior performance recently.

6.4.1 Modulation-Constrained Clustering Approach

In [15], a clustering-based LB classifier is proposed for modulation classification in multiple-input and multiple-output (MIMO) communication systems. In that work, a spatial-multiplexed MIMO system with N_t transmit antennas and N_r receive antennas is considered, in which data symbols are transmitted independently from each transmit antenna. The signal model of the n -th received signal vector $\mathbf{y}(n)$ is given as

$$\mathbf{y}(n) = \mathbf{H}\mathbf{s}(n) + \mathbf{u}(n), \quad n = 1, \dots, N \quad (6.15)$$

where $\mathbf{H} \in \mathbb{C}^{N_r \times N_t}$ is the channel matrix which remains constant within each block of N symbols, and $\mathbf{u}(n)$ denotes the AWGN vector.

For LB classifiers, the classification decision is made by selecting the modulation scheme with the maximum likelihood

$$\hat{M} = \arg \max_{M \in \mathcal{M}} \mathcal{L}_M \quad (6.16)$$

where \mathcal{L}_M is the likelihood function corresponding to the modulation scheme M and \mathcal{M} is the set of candidate modulation schemes.

Since the noise at the receiver is Gaussian, the PDF of the received signals follows the GMM given in (6.7) where $K = Q^{N_t}$, the mean and the Covariance matrix of the k -th Gaussian component are given as $\boldsymbol{\mu}_k = \mathbf{H}\mathbf{s}(n)$ and $\boldsymbol{\Sigma}_k = \sigma^2 \mathbf{I}$, respectively. The likelihood function for each modulation scheme can be calculated by estimating the parameters of the GMM model using the EM algorithm introduced in Sect. 6.2.1. However, the direct application of the EM algorithm presents the following challenges. Firstly, the modulation order Q of a modulation scheme determines the number of Gaussian components as well as the number of parameters to be estimated in the GMM model. Thus, the computational complexity for calculating the likelihood function of a higher-order modulation scheme can be extremely high. Secondly, the initialization of the set of parameters is an important part in the EM algorithm, and it would influence the converged performance and the convergence speed of the algorithm significantly. Hence, in order to improve the performance of the EM algorithm for modulation classification, there is a need to propose an EM algorithm with less parameters to be estimated and a good initialization method.

To reduce the number of parameters, a centroid reconstruction method is proposed in [15] by exploiting the relationship among the constellations. With the help of the proposed centroid reconstruction method, the number of parameters to be estimated is reduced from Q^{N_t} to N_t only. This also reduces the number of signal samples needed for the estimation. Specifically, for multiple-input and single-output (MISO) channels, the cluster centroids $\boldsymbol{\mu} = [\mu_1, \mu_2, \dots, \mu_K]$ can be reconstructed as follows

$$\boldsymbol{\mu} = \boldsymbol{\Psi} \mathbf{A} \quad (6.17)$$

where $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_{N_t}]^T$ is the reconstructive coefficient matrix, which is a known constant matrix for each modulation scheme, and $\boldsymbol{\Psi} = [r_1, \dots, r_{N_t}]$ is the corresponding reconstructive parameter vector.

By introducing constellation-structure-based centroid reconstruction in the EM algorithm, the iteration of $\{\boldsymbol{\mu}_k\}_{k=1}^K$ can be replaced by the iteration of r_1, r_2, \dots, r_{N_t} . If we denote $\boldsymbol{\Phi} = \{r_1, r_2, \dots, r_{N_t}, \delta^2 \mathbf{I}\}$ as the set of the unknown parameters, the likelihood function is shown as

$$\mathcal{L}(\boldsymbol{\Phi}) = \sum_{n=1}^N \ln \left(\sum_{k=1}^K \frac{1}{K} p(\mathbf{y}(n) | \boldsymbol{\Phi}) \right) \quad (6.18)$$

Hence, the proposed EM algorithm for modulation classification is shown as below.

1. **Initialization:** A two-stage initialization is proposed. In the first stage, fuzzy estimation is introduced to transform the parameters into a smaller range. In the second stage, modulation constrained K-means (MC K-means) algorithm is adopted, in which the constellation structure and the centroid reconstruction method are utilized.
2. **E-step:** Calculate the latent variable

$$\gamma_{nk} = \frac{p(\mathbf{y}(n) | \boldsymbol{\Phi})}{\sum_{k=1}^K \pi_k \cdot p(\mathbf{y}(n) | \boldsymbol{\Phi})} \quad (6.19)$$

and then the reconstructive parameters

$$r_1 = \frac{\sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} a_{1,k} (\mathbf{y}(n) - a_{2,k} r_2 - \dots - a_{N_t,k} r_{N_t})}{\sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} (a_{1,k})^2} \quad (6.20)$$

$$r_m = \frac{\sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} a_{m,k} (\mathbf{y}(n) - a_{1,k} r_1 - \cdots - a_{N_t,k} r_{N_t})}{\sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} (a_{m,k})^2} \quad (6.21)$$

where $m = 2, \dots, N_t$.

3. **M-step:** The cluster centroids $\boldsymbol{\mu}$ and the noise variance σ are updated iteratively as below

$$\boldsymbol{\mu}_k = \mathbf{a}_{1,k} r_1 + \mathbf{a}_{2,k} r_2 + \cdots + \mathbf{a}_{N_t,k} r_{N_t}, k = 1, \dots, K \quad (6.22)$$

and

$$\sigma^2 = \frac{\sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} (\mathbf{y}(n) - \boldsymbol{\mu}_k) (\mathbf{y}(n) - \boldsymbol{\mu}_k)^H}{\sum_{n=1}^N \sum_{k=1}^K \gamma_{nk}} \quad (6.23)$$

where $k = 1, \dots, K$.

4. **Classification Decision:** Repeat Step 2 and Step 3 iteratively until the likelihood function is converged. Then make the classification decision according to the criterion defined in (6.16).

Simulation results in [15] show that the proposed algorithm performs well with short observation length in terms of classification accuracy. Additionally, the performance achieved by the proposed algorithm is close to that of the average likelihood ratio-test upper bound (ALRT-UB), which can be seen as the performance upper bound of any modulation classification algorithm.

6.4.2 Deep Learning Approach

The modulation classifier in Sect. 6.4.1 requires accurate knowledge of the channel model. In addition, the channel model may not be available in practice. As a powerful supervised learning framework, DL can also be applied in modulation classification. In [16], a low-complexity blind data-driven modulation classifier based on DNN is proposed, which operates under uncertain noise condition modeled by a mixture of white Gaussian noise, white non-Gaussian noise and time-correlated non-Gaussian noise.

In [16], a single-input and single-output (SISO) channel is considered, and the n -th received signal sample is given as

$$r(n) = hs(n) + u(n), n = 1, 2, \dots, N \quad (6.24)$$

where $s(n)$ is the transmitted symbol from an unknown modulation scheme M_i , N is the number of symbols in a block, h is the channel coefficient and $u(n)$ denotes the additive noise.

Denote the set of the candidate modulation schemes and the received signal sequence by $\mathcal{M} = \{M_i, i = 1, 2, \dots, L\}$ and $\mathbf{r} = [r(1), r(2), \dots, r(N)]$, respectively. Let $P(M_i|\mathbf{r})$ denote the *a posteriori* probability of the modulation scheme M_i given the received signal \mathbf{r} . The objective of the work is to find the modulation scheme which maximizes the *a posteriori* probability. This is known as the maximum *a posteriori* (MAP) criterion.

$$\hat{M}_i = \arg \max_{M_i \in \mathcal{M}} P(M_i|\mathbf{r}) \quad (6.25)$$

In order to accurately make classification decisions with low complexity, the DNN is adopted to learn the *a posteriori* probability $P(M_i|\mathbf{r})$, $i = 1, \dots, L$. The DNN is used as an approximation function f mapping the received signal to the *a posteriori* probability. The in-phase and quadrature (IQ) components of the received signal samples are chosen as the inputs to the proposed neural network. Motivated by its superior performance for processing time-dependent data, the long short-term memory (LSTM) network is introduced in the design of the proposed neural network. There are three main reasons that the LSTM network is suitable for solving a modulation classification problem.

1. The LSTM network is able to learn features effectively from highly time-dependent data. This indicates the neural network with LSTM layer have advantages in learning the *a posteriori* probability from the signal samples which are highly time-dependent over time-correlated non-Gaussian channels.
2. Different from the fully-connected network which can only receive a one-dimensional as input, LSTM network allows two-dimensional vectors as the input in each time step. Hence, the network can process complex signal samples composed of IQ components, and can learn better from the input data.
3. Compared to the conventional fully-connected network, there are fewer parameters in the LSTM network because all time steps share the same weight matrices and biases.

Additionally, in order to summarize the output from the LSTM network, a temporal attention mechanism is adopted in the final LSTM layer over the outputs from all time steps. In the temporal attention mechanism, each output has a different weight, which indicates the importance of each to the modulation classification results.

Specifically, the proposed seven layer-neural network is composed of three stacked-LSTM layers and four fully-connected layers. In the training phase, the one-hot coding vectors of true modulation schemes of the input signal samples are used as the labels. The Adaptive Moment Estimation (Adam) optimizer is used to minimize the loss function to optimize the weights and bias in the network. After the training phase, the modulation classification is made by according to the MAP criterion defined in (6.25). The simulation results show that the classification accuracy

of the proposed classifier approaches that of the ML classifier with all the channel and noise parameters known. Moreover, under uncertain noise conditions, with lower computational online complexity, the proposed classifier can achieve a better performance than the EM and ECM classifiers.

6.5 Deep Reinforcement Learning for Dynamic Spectrum Access

In traditional DSA mechanism, there exists a centralized control node responsible for allocating the spectrum resources to users. Before making the access decisions, the centralized node needs to collect the global network information, such as the position information of users and base stations as well as the channel state information. However, such global network information is difficult to obtain in practice, as it imposes significant signal overheads on the system especially when there is a large number of users. Additionally, the collected information may be outdated in a highly dynamic network environment, resulting in invalid access strategy and poor performance. To solve the above issues, intelligent DSA framework operating with local network information is desirable. Recently, researchers introduced DRL techniques for DSA, showing superior performance on sequential decision-making tasks, to enable more flexible and intelligent DSA mechanism [17]. Since agents in DRL can make full use of the representation ability of neural networks, the decision space can be high-dimensional and continuous, which can guarantee the performance of the DSA mechanisms for large-scale networks.

In the following sections, we will introduce several typical applications on the use of DRL techniques for DSA.

6.5.1 Deep Multi-user Reinforcement Learning for Distributed Dynamic Spectrum Access

In [18], a DRL-based DSA framework is proposed to manage dynamic spectrum access in multichannel wireless networks, in which each user acts as an agent to make channel access decisions intelligently and independently to maximize its long-term transmission rate.

In this work, a wireless network composed of N users and K shared orthogonal channels is considered. Denote the set of users and the set of channels as $\mathcal{N} = \{1, 2, \dots, N\}$ and $\mathcal{K} = \{1, 2, \dots, K\}$, respectively. It is assumed that each user needs to choose a single channel for transmission in each time slot, and it always has packets to transmit. Additionally, the transmission is successful if there is only one user accessing the channel, and the transmission fails otherwise. After each transmission, each user can receive a binary observation $o_n(t)$ to indicate whether its transmission

is successful or not, i.e., $o_n(t) = 1$ if the transmission is successful and $o_n(t) = 0$ otherwise.

With the assumption that users don't have message exchange in each time slot, they can only make access decisions by their local observations. In order to solve the above problem, a DRL-based distributed framework for DSA is proposed, in which each user acts as an agent and constructs a DQN. The action space, state space and reward function are described as follows.

1. **Action Space:** In each time slot, each user needs to choose whether to transmit or not. If the user chooses to transmit, it needs to select a channel for transmission. The action of user n in time slot t is given as

$$a_n(t) \in \{0, 1, \dots, K\} \quad (6.26)$$

where $a_n(t) = 0$ indicates that user n chooses not to transmit in time slot t .

2. **State Space:** The state of each user is composed of its action and observation up to time slot t , which is given as

$$\mathcal{H}_n(t) = (\{a_n(i)\}_{i=1}^{t-1}, \{o_n(i)\}_{i=1}^{t-1}) \quad (6.27)$$

3. **Reward Function:** Since the objective is to maximize the long-term rate, the function of achievable rate is chosen as the reward function

$$r_n(t) = B \log_2(1 + SNR_n(k)) \quad (6.28)$$

where B is the channel bandwidth and $SNR_n(k)$ is SNR of user n on channel k .

In the DRL-based framework proposed in [18], in order to capture features from observations, the LSTM network is introduced in the structure of the adopted DQN. Additionally, the DDQN method is also adopted to improve the performance of the DQN. In the training phase, each user trains the parameters of their respective DQN cooperatively by communicating with a central unit. After updating the parameters, each user uses the trained DQN to make access decisions autonomously and independently. After the DQN is well-trained, the central unit is closed, and users use the converged DQN to obtain efficient access policy directly.

6.5.2 *Deep Reinforcement Learning for Joint User Association and Resource Allocation*

In heterogeneous networks (HetNets), all the base stations (BSs) normally provide services to users on shared spectrum bands in order to improve the spectrum efficiency. However, most existing methods need accurate global network information, e.g., channel state information, as the prior knowledge, which is difficult to obtain in practice.

In [19], a distributed DRL-based DSA framework is proposed for user association and resource allocation in the downlink HetNets. Specifically, a three-tier heterogeneous network is considered, which consists of N_m macrocell base stations (MBSs), N_p pico base stations (PBSs), N_f femto base stations (FBSs) and N user equipments (UEs). The sets of UEs and BSs are denoted, respectively, by $\mathcal{N} = \{1, \dots, N\}$ and $\mathcal{B} = \{0, 1, \dots, L-1\}$, where $L = N_m + N_p + N_f$. All the BSs share the same K orthogonal channels for downlink transmission, and the set of channels can be denoted as $\mathcal{K} = \{1, \dots, K\}$.

For each UE i , denote $b_i^l(t) = (b_i^0(t), \dots, b_i^{L-1}(t))$, $i \in \mathcal{N}$, $l \in \mathcal{B}$ as the binary *user-association* vector, where $b_i^l(t) = 1$ if UE i is associated with the BS l at time t and $b_i^l(t) = 0$ otherwise. For each BS, a binary *channel-allocation* vector is defined as $c_i^k(t) = (c_i^1(t), \dots, c_i^K(t))$, $i \in \mathcal{N}$, $k \in \mathcal{K}$, where $c_i^k(t) = 1$ if UE i uses channel resource C_k at time t and $c_i^k(t) = 0$ otherwise. It is assumed that each UE can only be connected to one BS and each channel can only be allocated to one UE for each BS in each time slot t .

The transmission power between UE i and its associated BS l on channel C_k at time t can be denoted as $p_{li}^k(t) = (p_{li}^1(t), \dots, p_{li}^K(t))$, $l \in \mathcal{B}$, $i \in \mathcal{N}$, $k \in \mathcal{K}$. Since all the BSs share the common spectrum resource, the co-channel interference should be considered. Hence, the signal-to-interference-plus-noise-ratio (SINR) of UE i associated with BS l and allocated with channel C_k is given as

$$\Gamma_{li}^k(t) = \frac{b_i^l(t) h_l^{i,k}(t) c_i^k(t) p_{li}^k(t)}{\sum_{j \in \mathcal{B} \setminus \{l\}} b_i^j(t) h_j^{i,k}(t) c_i^k(t) p_{ji}^k(t) + W N_0} \quad (6.29)$$

where $h_l^{i,k}(t)$ is the channel gain between the UE i and BS l at time t , W is the bandwidth of each channel and N_0 is the noise spectral power. Therefore, the total achievable transmission rate of UE i at time t can be expressed as

$$r_i(t) = \sum_{l=0}^{L-1} b_i^l(t) \sum_{k=1}^K W \log_2(1 + \Gamma_{li}^k(t)) \quad (6.30)$$

Considering that the operation cost of the UE i from BS l is determined by the transmit power $p_{li}^k(t)$, the total operation cost of UE i is given as

$$\varphi_i(t) = \sum_{l=0}^{L-1} \varphi_i^l(t) = \sum_{l=0}^{L-1} \lambda_l b_i^l(t) \sum_{k=1}^K c_i^k(t) p_{li}^k(t) \quad (6.31)$$

where λ_l is the price per unit of transmit power from BS l . Then we define the utility function of UE i as the total achievable profit minus the operation cost, which is denoted as

$$\omega_i(t) = \rho_i(t) r_i(t) - \varphi_i(t) \quad (6.32)$$

where $\rho_i > 0$ is the profit per unit transmission rate.

In this work, the objective of each UE is to maximize its own long-term utility. Since the problem is an integer programming problem and the objective of the problem is long-term, it is difficult to adopt the traditional optimization algorithms such as convex optimization to solve it. Additionally, the dimension of the decision space increases exponentially. Hence, a distributed DRL-based multi-agent framework for user association and resource allocation is proposed to maximize the long-term utility.

The state space, action space and reward function for modeling such a problem are given as follows.

1. **State Space:** In each time slot, the state is composed of the QoS of all the UEs, and we have

$$s(t) = \{s_1(t), s_2(t), \dots, s_N(t)\} \quad (6.33)$$

where $s_i(t)$ is a binary index indicating that whether UE i 's QoS is larger than the minimum threshold Ω_i or not, i.e., if the UE i 's QoS is larger than Ω_i , $s_i(t) = 1$ and otherwise, $s_i(t) = 0$.

2. **Action Space:** In each time t , each UE needs to choose a BS and a channel to access. Hence, the action of UE i consists of two parts, i.e, the user-association vector and the resource-allocation vector

$$a_{i_i}^k(t) = \{b_i^l(t), c_i^k(t)\} \quad (6.34)$$

where $b_i^l(t) \in \{0, 1\}$ and $c_i^k(t) \in \{0, 1\}$.

3. **Reward Function:** The reward function of UE i is mainly determined by its achievable rate in time slot t . Besides, to improve the convergence performance of the algorithm, the action-selection cost is also considered in the design of the reward function.

$$R_i(t) = \begin{cases} \omega_i(t), & \Gamma_i(t) \geq \Omega_i \\ -\Psi_i, & \text{otherwise} \end{cases} \quad (6.35)$$

where $\Gamma_i = \sum_{l=0}^{L-1} \sum_{k=1}^K \Gamma_k^{l,i}$ is SINR of UE i , Ω_i is a pre-designed minimum QoS requirement and Ψ_i is the action-selection cost, which is a positive value.

In the proposed framework, each UE is equipped with a DQN to make access decisions independently. In the initialization stage, each UE is first connected to the BS which resulted in the maximum received signal reference power (RSRP) and constructs a DQN, in which the parameter is initialized randomly. At each training time t , each UE has a common state s and selects an action, namely, access request, according to its Q-value $Q_i(s, a_i, \theta)$ obtained from its DDQN. The access request contains the indices of the required BS and the channel. Then, if the BS accepts the request, the BS would send a feedback signal to the UE, which indicates the resource is available, and otherwise, the BS would not reply. After connecting to the chosen BS and accessing the chosen channel, the UE obtains an immediate reward $u_i(s, a_i)$ and a new state s' , then stores the current experience $\langle s, a_i, u_i(s, a_i), s' \rangle$

into its replay memory \mathcal{D} . Finally, each UE updates the parameter θ of its DDQN by using stochastic gradient descent (SGD) algorithms based on the random samples from the memory \mathcal{D} .

6.6 Summary

In this chapter, we have provided a brief review on machine learning techniques, and have described some applications on AI-based DSM mechanisms such as spectrum sensing, signal classification and dynamic spectrum access. These AI-based DSM mechanism have been shown to achieve better performance and robustness than conventional schemes. Additionally, they can also provide more efficient and flexible ways to implement the DSM. In the future, the combination of AI techniques and the DSM mechanisms would become a novel and promising research direction.

References

1. S. Russell, N. Peter, *Artificial Intelligence: A Modern Approach*, 3rd edn. (Prentice Hall Press, Upper Saddle River, 2009)
2. What is artificial intelligence, <https://www.aisb.org.uk/public-engagement/what-is-ai>
3. The 5G future will be powered by AI, <https://www.networkcomputing.com/wireless-infrastructure/5g-future-will-be-powered-ai>
4. Leveraging machine learning and artificial intelligence for 5G, <https://www.cablelabs.com/leveraging-machine-learning-and-artificial-intelligence-for-5g>
5. C. Jiang, H. Zhang, Y. Ren, Z. Han, K. Chen, L. Hanzo, Machine learning paradigms for next-generation wireless networks. *IEEE Wirel. Commun.* **24**(2), 98–105 (2017)
6. T.M. Mitchell, *Machine Learning* (McGraw-Hill, New York, 1997)
7. R.S. Sutton, A.G. Barto, *Introduction to Reinforcement Learning*, 1st edn. (MIT Press, Cambridge, 1998)
8. C.M. Bishop, *Pattern Recognition and Machine Learning* (Springer, New York, 2006)
9. Y. LeCun, Y. Bengio, G. Hinton, Deep learning. *Nature* (2015)
10. A. Krizhevsky, I. Sutskever, G. Hinton, Imagenet classification with deep convolutional neural networks, in *Proceedings of the NIPS* (2012)
11. A. Graves, A. Mohamed, G. Hinton, Speech recognition with deep recurrent neural networks, in *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), BC, Vancouver* (2013), pp. 6645–6649
12. V. Mnih et al., Human-level control through deep reinforcement learning. *Nature* **518**(7540), 529–533 (2015)
13. H. Van Hasselt, A. Guez, D. Silver, Deep reinforcement learning with double Q-learning, in *Proceedings of the AAAI* (2016), pp. 2094–2100
14. K.M. Thilina, K.W. Choi, N. Saquib, E. Hossain, Machine learning techniques for cooperative spectrum sensing in cognitive radio networks. *IEEE J. Sel. Areas Commun.* **31**(11), 2209–2221 (2013)
15. J. Tian, Y. Pei, Y. Huang, Y.-C. Liang, Modulation-constrained clustering approach to blind modulation classification for MIMO systems. *IEEE Trans. Cogn. Commun. Netw.* **4**(4), 894–907 (2018)

16. S. Hu, Y. Pei, P.P. Liang, Y.-C. Liang, Robust modulation classification under uncertain noise condition using recurrent neural network, in *Proceedings of the IEEE Global Communications Conference (GLOBECOM'18), United Arab Emirates, Abu Dhabi* (2018), pp. 1–7
17. N.C. Luong, D.T. Hoang, S. Gong, D. Niyato, P. Wang, Y.-C. Liang, D.I. Kim, Applications of deep reinforcement learning in communications and networking: a survey. *IEEE Commun. Surv. Tutor.* (2019)
18. O. Naparstek, K. Cohen, Deep multi-user reinforcement learning for distributed dynamic spectrum access. *IEEE Trans. Wirel. Commun.* **18**(1), 310–323 (2019)
19. N. Zhao, Y.-C. Liang, D. Niyato, Y. Pei, Y. Jiang, Deep reinforcement learning for user association and resource allocation in heterogeneous networks, in *Proceedings of the IEEE Global Communications Conference (GLOBECOM'18), United Arab Emirates, Abu Dhabi* (2018), pp. 1–6

Open Access This chapter is licensed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

