

Verstärkendes Lernen

Unter einem sequenziellen Entscheidungsprozess (engl. SDP) versteht man jegliche Art von dynamischen Prozessen, bei denen in einer festgelegten zeitlichen Reihenfolge mehrere aufeinanderfolgende Entscheidungen getroffen werden, aus denen sich wiederum Bestrafungen bzw. Belohnungen ergeben, die sich mit der Zeit anhäufen. Es geht darum, die Entscheidungssequenz zu bestimmen, welche die definierten Belohnungen maximiert.

Eine dritte Art des maschinellen Lernens, bekannt als verstärkendes Lernen, überschneidet sich in einigen grundlegenden Eigenschaften mit denen des überwachten und des unüberwachten Lernens, jedoch wird sie gezielt für sequenzielle Entscheidungsprozesse eingesetzt, da hier auf andere Arten von Daten zurückgegriffen wird als beim überwachten und unüberwachten Lernen. Ein Großteil der Theorie des verstärkenden Lernens beruht auf dem Prinzip der **dynamischen Programmierung**, das in den 1950er Jahren von Richard Bellman eingeführt wurde (Bellman, 1957). Die dynamische Programmierung bildet zwar, wie wir noch sehen werden, die mathematische Grundlage für das verstärkende Lernen, allerdings ist ihre Rechenkapazität beschränkt. In der Praxis bedeutet dies, dass die dynamische Programmierung, zumindest in ihrer Reinform, keine gangbare Lösung für Anwendungen bietet, die den heutzutage üblichen hohen Datendurchsatz nutzen.

Eine besonders überzeugende frühe Anwendung des verstärkenden Lernens war die Entwicklung eines SDPs namens TD-Gammon. Dieser ist in der Lage, Backgammon auf wettbewerbsfähigem Niveau zu spielen. In der jüngeren Vergangenheit verwendete DeepMind Technologies die Methode des verstärkenden Lernens, um ein Schachprogramm zu entwickeln, das unter dem Namen AlphaZero bekannt wurde. Das verstärkende Lernen ist besonders gut für die Spieleentwicklung geeignet, vor allem in Verbindung mit tiefen neuronalen Netzwerkmodellen (tiefes verstärkendes Lernen). Derzeit wird es beispielsweise in der Entwicklung von sehr effektiven Videospiel-Apps genutzt.