

Markov
De-
ci-
sion
Pro-
cesses
We

de-
fine
a
con-
trolled
sys-
tem
as
a
la-
beled
Markov
de-
ci-
sion
pro-
cess.

[Labeled
Markov
De-
ci-
sion
Pro-
cess]

A
(la-
beled)
Markov
de-
ci-
sion
pro-
cess
(MDP)
is
a
tu-
ple
 M

=
 $(S, A, \mathcal{A}, P, s_{init}, AP, L),$
where

S
is
a
fi-
nite
set
of
states,

A
is
a
fi-
nite
set
of
ac-
tions,

$\mathcal{A} : S \rightarrow 2^A$

is
a
map-
ping
that
maps
each
state

a-
 bil-
 ity
 such
 that

$$\sum_{s' \in S} P(s'|s, a) = 1$$
 for
 any
 state
 $s \in S$
 and
 any
 ac-
 tion
 $a \in \mathcal{A}(s)$,
 $s_{init} \in S$
 is
 the
 ini-
 tial
 state,
 AP
 is
 a
 fi-
 nite
 set
 of
 atomic
 propo-
 si-
 tions,
 and
 $L : S \times A \times S \rightarrow 2^{AP}$
 is
 a
 la-
 bel-
 ing
 func-
 tion
 that
 as-
 signs
 a
 set
 of
 atomic
 propo-
 si-
 tions
 to
 each
 tran-
 si-
 tion
 $(s, a, s') \in S \times A \times S$.

In
 the
 MDP
 M ,
 an
 in-
 fi-
 nite
 path

action-
value
func-
tion
 $Q^\pi(s, a)$
un-
der
the
pol-
icy
 π
as
fol-
lows.

$$Q^\pi(s, a) = E^\pi[\sum_{n=0}^{\infty} \gamma^n \mathcal{R}(S_n, A_n, S_{n+1}) | S_0 = s, A_0 = a].$$

We
have
the
fol-
low-
ing
re-
cur-
sively
equa-
tion
for
the
state-
value
func-
tion
and
the
action-
value
func-
tion.

$$\begin{aligned} V^\pi(s) &= E^\pi[\sum_{n=0}^{\infty} \gamma^n \mathcal{R}(S_n, A_n, S_{n+1}) | S_0 = s] \\ &= \sum_{a \in \mathcal{A}(s)} \pi(s, a) \sum_{s' \in S} P(s' | s, a) E^\pi[\sum_{n=0}^{\infty} \gamma^n \mathcal{R}(S_n, A_n, S_{n+1}) | S_0 = s, A_0 = a, S_1 = s'] \\ &= \sum_{a \in \mathcal{A}(s)} \pi(s, a) \sum_{s' \in S} P(s' | s, a) \{ \mathcal{R}(s, a, s') + \gamma E^\pi[\sum_{n=0}^{\infty} \gamma^n \mathcal{R}(S_n, A_n, S_{n+1}) | S_1 = s'] \} \\ &= \sum_{a \in \mathcal{A}(s)} \pi(s, a) \sum_{s' \in S} P(s' | s, a) \{ \mathcal{R}(s, a, s') + \gamma V^\pi(s') \}, \end{aligned}$$

by
the
def-
i-
ni-
tion
of
the
action-
value
func-
tion,
it
holds
that

$$\begin{aligned} Q^\pi(s, a) &= \max_{a \in \mathcal{A}(s)} V^\pi(s) \\ &= \sum_{s' \in S} P(s' | s, a) \{ \mathcal{R}(s, a, s') + \end{aligned}$$

over-
all
pro-
ce-
dure
TD-
learning
for
a
state-
value
func-
tion
is
given
by
Al-
go-
r-
i-
tym
.

TD-
learning
meth-
ods
for
an
action-
value
func-
tion
are
clas-
si-
fied
as
two
main
learn-
ing
meth-
ods
that
are
re-
ferred
to
Q-
learning
and
SARSA.

Stochastic
Dis-
crete
Event
Sys-
tems
We
rep-
re-
sent
a
stochas-
tic
dis-
crete
event
sys-
tem
(DES)
as
an
MDP.

[Stochastic
dis-
crete

a-
tive
fre-
quency
of
oc-
cur-
rence
of
each
event
does
not
de-
pend
on
the
con-
trol
pat-
tern.

We
de-
fine
a
re-
ward
func-
tion
 $\mathcal{R} :$
 $S \times$
 $2^E \times$
 $E \times$
 $S \rightarrow$
 R
and
the
re-
ward
 \mathcal{R}
can
be
de-
com-
posed
into
 \mathcal{R}_1
and
 \mathcal{R}_2 .
The
first
re-
ward
 $\mathcal{R}_1 :$
 $S \times$
 $2^E \rightarrow$
 R
is
de-
ter-
mined
by
the
con-
trol
pat-
tern
se-
lected
by
the
su-
per-
vi-
sor,
which