**Soltan A. ALIEV, Pavlo R. Shpak, Yaroslav I. YELEYKO**

# ASSESSMENT AND OPTIMAL POLICIES OF LIMIT KILLED MARKOV DECISION PROCESS

### Abstract

*We analyzed limit killed Markov decision process on finite time interval for finite and countable models, as well as its assessment and optimal policy. We showed that assessment of limit process is equal to limit of assessments, both for countable and finite model. For finite models we proved that limit point of sequence of optimal policies always exists and each limit point is optimal policy for limit process. For countable models we proved that if limit point of $\varepsilon$ -optimal policies exists than this point is $\delta$ - optimal policy for limit process, where $\delta$ is arbitrary real number greater than $\varepsilon$.*

### 1. Introduction

Bellman's book [6] can be considered as a starting point of Markov decision processes, where ideas of dynamic programming principle were applied to Markov decision processes. Later Markov decision processes are well described in [1]: the definition of Markov decision process was given, as well as definition of optimal policy, assessment of the policy and assessment of process. But in [1] the model does not take into account the risk factor, namely the probability of bankruptcy at some determined moment of time. As a result, we come to the idea of killed Markov decision process where the business can crash with some nonzero probability at every moment of time, with the exception of the initial state. The basic ideas about killing of Markov processes is given in [5]. Also some aspects of Markov decision processes where described by Feinberg and Shwartz [4].

### 2. Definitions

Let $X_t(t=m,\ldots,n)$ and let $A_t(t=m+1,\ldots,n)$ be countable or finite.

**Remark.** All definitions and basic ideas of killed Markov decision process are given accordingly to [2] and [3].

**Definition.** *A mapping $T: X_t \times A_{t+1} \times X_{t+1} \to P(X_t)$ which satisfies following conditions:*

$$T(x,a,x') = \mathbb{P}(x_{t+1} = x \mid x_t = x', a_{t+1} = a) \equiv \mathbb{P}$$

*is called **transition function** and we we'll it as $p(x|a)$.*

**Definition.** *The point $x^* \in X_t$ is called **killed state**, and $p(x^*|a)$ - **probability of kill***

**Definition.** *The trajectory $l = x_m a_{m+1} x_{m+1} \ldots a_t x_t$ $n$ is called way if $t = n$ or $x_t = x^*$. The set of all ways we'll denote as $L = X \times (A \times X)^{n-m}$.*

**Definition [Killed Markov decision process].** *A killed Markov decision process on a time interval $[m,n]$ is defined through the following objects:*

1. *Sets $X = \bigcup\limits_{t=m}^{n} X_t$ (spaces of states);*
2. *Sets $A = \bigcup\limits_{t=m+1}^{n} A_t$ (spaces of actions);*
3. *The projection mapping $j: A \to X$, where $j(A_{t+1}) = X_t \backslash \{x^*\}, x^* \in X_t$,*
4. *Probability distribution $p(\cdot|a) = \mathbb{P}(x_t = x|a_t = a, x_{t-1})$ on $X_t$ with killed states:*

$$\mathbb{P}(x_{t+1} = x^* \mid a_t = a) \equiv p(x^* \mid a) \geq 0; ;$$

5. *Function $q: A \to \mathbb{R}$ (reward function);*

*6. Function $r : X_n \to \mathbb{R}$ (terminal reward);*

*7. Function c (crash function), defined on killed states:*

$$c(x) = \begin{cases} -\sum\limits_{i=m+1}^{t} \sup\limits_{a_i \in A_i} q(a_i) & , \ x = x^* \\ 0 & , \ x \neq x^* \end{cases}$$

$x \in X_t$, $x^*$ - *killed state. (function c ensures a total bankruptcy - total loss of accumulated capital or more);*

*8. Initial distribution $\mu$ on $X_m$.*

*A stochastic process defined through (1-8) is called* **killed Markov decision process** *and is denoted as $Z_\mu$. If the initial distribution $\mu$ is concentrated in the point x, we shall write $Z_x$.*

**Definition.** *Let $A(x) \subset A$ - is the set of all available actions at state $x \in X$. $\varphi(x) : X \backslash x^* \to A(x)$ is called* **simple policy** *if $\varphi(x_{t-1}) = a_t$ $\forall t = m+1, \ldots, n$.*

**Definition.** *The mapping $\pi : H \to \pi(\cdot | h \in H)$ is called* **policy**, *where $\pi(\cdot | h \in H)$ - probability distribution on $A(x_{t-1})$ and H - the space of histories up to epoch $(h \in H \Leftrightarrow h = x_m a_{m+1}, \ldots, a_{t-1} x_{t-1})$.*

**Definition.** *Policy $\pi(\cdot | h)$ is called* **Markov policy** *if $\pi(\cdot | h) = \pi(\cdot | x_{t-1})$.*

**Assumption.** *The reward function q and terminal reward function r have the supremum, $\exists \sup_{a \in A} q(a)$ and $\exists \sup_{x \in X_n} r(x)$.*

If $p(\cdot | a)$ is the transition function and $\pi(\cdot | h)$ is a policy, then for $\forall \mu$ - initial distribution exists corresponding probability distribution $P$ defined on space $L$ which has such notation:

$$P(l) = \begin{cases} \mu(x_m) \pi(a_{m+1} | x_m) p(x_{m+1} | a_{m+1}) \cdots \\ \quad \cdot \pi(a_n | h_{n-1}) p(x_n | a_n), & , \ \forall t \ \ x_t \neq x^* \\ \mu(x_m) \pi(a_{m+1} | x_m) p(x_{m+1} | a_{m+1}) \cdots \\ \quad \cdot \pi(a_t | h_{t-1}) p(x_t^* | a_t), & , \ \ x_t = x^* \end{cases}$$

For all function $\xi$ from space $L$ the mathematical expectation of $\xi$ is:

$$E^*(\xi) = \sum_{l \in L} \xi(l) \, \mathrm{P}(l) \quad (2.1.2)$$

The assessment

$$I(l) = \begin{cases} \sum\limits_{t=m+1}^{n} q(a_t) + r(x_n) & , \forall r \ \ x_k \neq x^* \\ \sum\limits_{k=m+1}^{t} q(a_t) + c(x_t^*) & , \ \ x_t = x^* \end{cases}$$

of the way $l$ is example of such function. And we'll denote its expectation as $\omega$:

$$\omega = \mathrm{E} I(l)$$

**Definition [Assessment of policy].** *The value $\omega$ is called* **assessment of policy** .

**Definition [Assessment of process].** *$\upsilon \equiv \sup_\pi \omega(\pi)$ is called* **assessment of killed Markov decision process** *$Z_\mu^*$ or* **assessment of initial distribution** *$\mu$.*

**Definition [optimal policy].** *If spaces of states and spaces of actions are finite sets then policy $\pi$ is called* **optimal** *if $\omega(\pi) = \upsilon$.*

**Definition [$\varepsilon$-optimal policy].** *If spaces of states and spaces of actions are countable or finite sets and at least one of them is countable then policy $\pi$ is called $\varepsilon$-optimal if $\omega(\pi) \geq \upsilon - \varepsilon$.*

### 3. Formulation of the Problem

Lets consider that we have a family of killed Markov decision processes $\left(Z_\mu^k\right)_{k=0}^\infty$ on a time interval $[m, n]$, where $Z_\mu^k = (X, A, j, p_k, q_k, r_k, c_k, \mu_k)$, and following conditions are hold:

$$\forall t, \ \forall a_t \in A_t, \ \forall x_t \in X_t : \lim_{k \to \infty} p_k \left( x_t \mid a_t \right) = p \left( x_t \mid a_t \right),$$

$$\forall t, \ \forall a_t \in A_t : \lim_{k \to \infty} q_k \left( a_t \right) = q \left( a_t \right),$$

$$\forall x_n \in X_n : \lim_{k \to \infty} r_k \left( x_n \right) = r \left( x_n \right),$$

$$\forall x_m \in X_m : \lim_{k \to \infty} \mu_k \left( x_m \right) = \mu \left( x_m \right),$$

$$\forall t : \lim_{k \to \infty} c_k \left( x_t^* \right) = c \left( x_t^* \right),$$

$$\forall k : \sum_{x_m \in X_m} \mu_k(x_m) = 1,$$

$$\forall t, \ \forall a_t \in A_t, \ \forall k : \sum_{x_t \in X_t} p_k \left( x_t \mid a_t \right) = 1.$$

Our goal is to find assessment and optimal ($\varepsilon$-optimal) policies of limit process $Z_\mu = (X, A, j, p, q, r, c, \mu)$.

### 4. Finite Case

Let $X_t (t=m, \ldots, n)$ and let $A_t (t=m+1, \ldots, n)$ be countable sets.

**Theorem 1.** *Let $\left(Z_\mu^k\right)_{k=0}^\infty$ be a family of killed Markov decision processes as defined above, then:*

*1. $\upsilon = \lim_{k \to \infty} \upsilon_k$ , where $\upsilon$ - assessment of limit process $Z_\mu$ and $\upsilon_k$ assessments of corresponding processes $Z_\mu^k$.*

*2. Exists at least one limit point of $(\overline{\pi}_k)_{k=0}^\infty$, where $\overline{\pi}_k$ – optimal policy for corresponding processes $Z_\mu^k$*

*3. If $\overline{\pi}$ is a limit point of $(\overline{\pi}_k)_{k=0}^\infty$ then $\overline{\pi}$ is optimal policy for limit process $Z_\mu$.*

**Proof.** For $\forall k$ exists optimal policy $\overline{\pi}_k$ of corresponding process $Z_\mu^k$ which satisfies following conditions ([3]):

$$\omega_k \left( \mu_k, \ \overline{\pi}_k \right) = \sup_\pi \omega_k \left( \mu_k, \ \pi \right) = \upsilon_k(\mu_k)$$

$$\forall t, \ \forall a_{t+1} \in A_{t+1}, \ \forall x_t \in X_t : \ 0 \leq \overline{\pi}_k \left( a_{t+1} \mid x_t \right) \leq 1.$$

Let's consider a sequence of optimal polices $(\overline{\pi}_k)_{k=1}^\infty$. As $X$ and $A$ are finite sets $(\overline{\pi}_k)_{k=1}^\infty$ is a subset of finite Cartesian product of compacts, thus exists convergent sub sequence $(\overline{\pi}_{k_s})_{s=1}^\infty$, which satisfies following conditions:

$$\forall t, \ \forall x_t \in X_t : \sum_{a_{t+1} \in A_t} \overline{\pi}_{k_s} \left( a_{t+1} \mid x_t \right) = 1 \quad ,$$

$$\forall t, \ \forall a_{t+1} \in A_{t+1}, \ \forall x_t \in X_t : \ \overline{\pi}_{k_s} \left( a_{t+1} \mid x_t \right) \xrightarrow[s \to \infty]{} \overline{\pi} \left( a_{t+1} \mid x_t \right)$$

For limit process $Z_\mu = (X, A, j, p, q, r, c, \mu)$ exists optimal policy $\overline{\pi}^*$:

$$\omega(\mu, \ \overline{\pi}^*) = \sup_\pi \omega(\mu, \ \pi)$$

We need to show that assessment of policy $\overline{\pi}$ is equal to assessment of limit process $Z_\mu$:

$$\omega(\mu, \overline{\pi}) = \omega(\mu, \overline{\pi}^*) = \upsilon(\mu)$$

Let's suppose that this is not true:

$$\omega(\mu, \overline{\pi}) \neq \upsilon(\mu)$$

Let $\omega(\mu, \overline{\pi}) > \upsilon(\mu)$. As $\upsilon(\mu) = \omega(\mu, \ \overline{\pi}^*) = \sup_\pi \omega(\mu, \ \pi) \geq \omega(\mu, \ \pi')$ for $\forall \pi'$, consequently $\omega(\mu, \overline{\pi}) \leq \upsilon(\mu)$.

Let $\omega(\mu, \overline{\pi}) < \upsilon(\mu)$ and $\omega(\mu, \ \overline{\pi}^*) - \omega(\mu, \overline{\pi}) = \ \varepsilon > 0$. As

$$\mu_k \underset{k \to \infty}{\longrightarrow} \mu, \ p_k \underset{k \to \infty}{\longrightarrow} p, \ q_k \underset{k \to \infty}{\longrightarrow} q, \ c_k \underset{k \to \infty}{\longrightarrow} c, \ r_k \underset{k \to \infty}{\longrightarrow} r, \ \overline{\pi}_{k_s} \underset{s \to \infty}{\longrightarrow} \overline{\pi}.$$

Then:

$$\forall \ \delta > 0 \ \exists \ k_0 \ \forall \ k > \ k_0 : \ \left|\omega_k(\mu_k, \ \overline{\pi}^*) - \omega(\mu, \ \overline{\pi}^*)\right| < \ \delta,$$

$$\forall \ \delta > 0 \ \exists \ s_0 \ \forall \ s > \ s_0 : \ \left|\omega_{k_s}(\mu_{k_s}, \ \overline{\pi}_{k_s}) - \omega(\mu, \ \overline{\pi})\right| < \ \delta.$$

And

$$\forall k : \ 0 < \omega(\mu, \overline{\pi}^*) - \omega(\mu, \overline{\pi}) = \ \omega(\mu, \overline{\pi}^*) - \omega_k(\mu, \overline{\pi}^*) +$$
$$+ \ \omega_k(\mu, \overline{\pi}^*) - \omega_k(\mu_k, \ \overline{\pi}_k) + \omega_k(\mu_k, \ \overline{\pi}_k) - \omega(\mu, \overline{\pi})$$

Moreover $\overline{\pi}_k$ is optimal policy for process $Z_\mu^k$, thus $\forall k : \ \omega_k(\mu, \overline{\pi}^*) - \omega_k(\mu_k, \ \overline{\pi}_k) \leq 0$, as a result:

$$\forall k \ : \ \ 0 < \omega(\mu, \overline{\pi}^*) - \omega(\mu, \overline{\pi}) \leq \ \omega(\mu, \overline{\pi}^*) - \omega_k(\mu, \overline{\pi}^*) + \omega_k(\mu_k, \ \overline{\pi}_k) - \omega(\mu, \overline{\pi}) \leq$$
$$\leq \ \left|\omega_k(\mu_k, \ \overline{\pi}^*) - \omega(\mu, \ \overline{\pi}^*)\right| + \ \left|\omega_k(\mu_k, \ \overline{\pi}_k) - \omega(\mu, \ \overline{\pi})\right|$$

$$\forall \ k \ > k^* = \max\{k_0, \ \ k_{s_o}\} \ : 2\delta > \left|\omega_k(\mu_k, \ \overline{\pi}^*) - \omega(\mu, \ \overline{\pi}^*)\right| +$$

$$+ \ \left|\omega_k(\mu_k, \ \overline{\pi}_k) - \omega_k(\mu_k, \ \overline{\pi})\right| \geq \left|\omega(\mu, \overline{\pi}^*) - \omega(\mu, \overline{\pi})\right| = \ \omega(\mu, \ \overline{\pi}^*) - \omega(\mu, \overline{\pi}) = \ \varepsilon > 0$$

Which means that $\omega(\mu, \overline{\pi}) = \omega(\mu, \overline{\pi}^*) = \upsilon(\mu)$. Theorem is proved.

### 5. Countable Case

Let $X_t(t{=}m, \dots, n)$ and let $A_t(t{=}m{+}1, \dots, n)$ be countable or finite sets and at least one of them is countable.

$\exists \ sup_{a \in A} \ q(a)$ and $\exists \ sup_{x \in X_n} \ r(x)$

**Theorem 2.** *Let $\left(Z_\mu^k\right)_{k=0}^\infty$ be a family of killed Markov decision processes as defined above. If $\forall k \ \exists \ sup_{a \in A} \ q_k(a) \neq \infty$, $\exists \ sup_{x \in X_n} r_k(x)$ and $\exists \ sup_{a \in A} \ q(a) \neq \infty$, $\exists \ sup_{x \in X_n} r(x)$ then:*

*1. $\upsilon = \ lim_{k \to \infty} \upsilon_k$, where $\upsilon$ - assessment of limit process $Z_\mu$ and $\upsilon_k$ assessments of corresponding processes $Z_\mu^k$.*

*2. If exists at limit point of $(\overline{\pi}_k)_{k=0}^\infty$, where $\overline{\pi}_k - \varepsilon$ - optimal policy for corresponding processes $Z_\mu^k$ then this limit point is $\delta$ - optimal policy for limit process $Z_\mu$, where $\delta$ is arbitrary real number greater than $\varepsilon$.*

**Proof.** As

$$v_k\left(\mu_k\right) = \sup_\pi \omega_k\left(\mu_k,\ \pi\right) = \sup_\pi \mathrm{E} I\left(l\right) = \sup_\pi \sum_{l\in L} I\left(l\right) P_k\left(l\right)$$

Let $L'$ be a set of way where $\forall t\ \ x_t \neq x^*$ and $L'' = L\backslash L'$ then

$$v_k\left(\mu_k\right) = \sup_\pi \left(\sum_{l\in L'} I\left(l\right)P_k\left(l\right) + \sum_{l\in L''} I\left(l\right)P_k\left(l\right)\right) =$$

$$= \sup_\pi \left(\sum_{l\in L'}\left(\sum_{i=m+1}^n q_k\left(a_i\right)+r_k\left(x_n\right)\right)\times\right.$$

$$\times\left(\mu_k\left(x_m\right)\pi\left(a_{m+1}\mid x_m\right)\ldots\pi\left(a_n\mid h_{n-1}\right)p_k\left(x_n\mid a_n\right)\right)+$$

$$+ \sum_{l\in L''}\left(\sum_{i=m+1}^t q_k\left(a_i\right)+c_k\left(x_t^*\right)\right)\times$$

$$\times\left(\mu_k\left(x_m\right)\pi\left(a_{m+1}\mid x_m\right)\ldots\pi\left(a_n\mid h_{n-1}\right)p_k\left(x_t^*\mid a_n\right)\right)\bigg)$$

As a result

$$\lim_{k\to\infty}v_k\left(\mu_k\right) = \sup_\pi \left(\sum_{l\in L'}\left(\sum_{i=m+1}^n q\left(a_i\right)+r\left(x_n\right)\right)\times\right.$$

$$\times\left(\mu\left(x_m\right)\pi\left(a_{m+1}\mid x_m\right)\ldots\pi\left(a_n\mid h_{n-1}\right)p\left(x_n\mid a_n\right)\right)+$$

$$+ \sum_{l\in L''}\left(\sum_{i=m+1}^t q\left(a_i\right)+c\left(x_t^*\right)\right)\left(\mu\left(x_m\right)\pi\left(a_{m+1}\mid x_m\right)\times\right.$$

$$\times\pi\left(a_n\mid h_{n-1}\right)p\left(x_t^*\mid a_n\right)\bigg)\bigg) = v(\mu)$$

In addition if exists at limit point of $\left(\overline{\pi}_k\right)_{k=0}^\infty$ this means that exists convergent subsequence $\left(\overline{\pi}_{k_s}\right)_{s=1}^\infty$, which satisfies following conditions:

$$\forall t,\ \forall x_t \in X_t:\ \sum_{a_{t+1}\ \in\ A_t} \overline{\pi}_{k_s}\left(a_{t+1}\mid x_t\right) = 1\ \ ,$$

$$\forall t,\ \forall a_{t+1} \in\ A_{t+1},\ \forall x_t \in X_t:\ \overline{\pi}_{k_s}\left(a_{t+1}\mid x_t\right)\xrightarrow[s\to\infty]{}\overline{\pi}\left(a_{t+1}\mid x_t\right)$$

$$\forall k_s:\ \omega_{k_s}\left(\mu_{k_s},\ \overline{\pi}_{k_s}\right) > v_{k_s}\left(\mu_{k_s}\right) - \varepsilon$$

So

$$\omega\left(\mu,\ \overline{\pi}\right) = \lim_{s\to\infty}\omega_{k_s}\left(\mu_{k_s},\ \overline{\pi}_{k_s}\right)\ \geq\ v\left(\mu\right) - \varepsilon$$

Consequently $\overline{\pi}$ is $\delta$ – optimal for limit process $Z_\mu$ where $\delta$ is arbitrary real number greater than $\varepsilon$. Theorem is proved.

## 6. Example
Let's consider a family of killed Markov decision processes $\left(Z_\mu^k\right)_{k=0}^\infty$ where $Z_\mu^k$:

From definition of crash function $c(x) = -\sum\limits_{i=m+1}^{t} \sup_{a_i \in A_i} q(a_i)$ for $x_t = x^*$ we receive that $c(x_1^*) = -15$ and $c(x_2^*) = -30 - (-1)^n 3/n$. As a result assessments of states $(III)$, $(IV)$ and $(V)$ would equal to following values:

$$v_n(III) = \max?\{22 + (-1)^n 12/5n, \ 22 + 1/2n - (-1)^n 1/2n\}$$

$$v_n(IV) = 16 - (-1)^n 3/10n$$

$$v_n(V) = 16 - (-1)^n 3/4n.$$

Accordingly $v_n(III) = 22 + 12/5n$ if $n$ is even and $v_n(III) = 22 - 1/n$ if $n$ is odd number.

Thus is $n$ is even assessments of states $(I)$ and $(II)$ would be as following:

$$v_n(I) = 13 + 41/40n + n/n + 1$$

$$v_n(II) = 24 + 193/32n + 9/20n^2.$$

And if $n$ is odd:

$$v_n(I) = 13 - 67/120n + n/n + 1$$

$$v_n(II) = 24 + 783/160n - 17/80n^2.$$

As $v_n(\mu_n) = \mu_n(I) v_n(I) + \mu_n(II) v_n(II)$

We receive that assessment of limit process equals to:

$$v(\mu) = 13\mu(I) + 24\mu(II).$$

### References

[1]. Dynkin E.B., Yushkevich A.A. *Markov Decision Processes*, M., 1975, 334 p.(Russian).

[2]. Parolya N. R., Yeleyko Y. I. *Killed Markov decision processes on finite time interval for countable models*, Transactions of NAS of Azerbaijan, 2010, vol. XXX, No 4, pp. 141-152.

[3] Parolya N. R., Yeleyko Y. I. *Killed Markov decision processes on finite time interval for finite models*, Visnyk of the Lviv Univ. Series Mech. Math. 2013. Issue 72, pp. 243-254. (Ukrainian)

[4]. Feinberg E.A., Shwartz A. *Handbook of Markov Decision Processes*, Kluwer, 2002, 565 p.

[5] Pakes A.G., *Killing and Resurrection of Markov Processes*, Stochastic Models, 1997, v.13, I.2, pp.255-269.

[6]. Bellman R.E. *Dynamic Programming*, Princeton University Press, 1957, 400 p.

**Soltan A. Aliev**
Institute of Mathematics and Mechanics of NAS of Azerbaijan
9, B. Vahabzade str., AZ1141, Baku, Azerbaijan
Tel: (99412) 438 22 44 (off.).

**Pavlo R. Shpak, Yaroslav I. Yeleyko**
Ivan Franko National University of Lviv
1, Universytetska str., 79000, Lviv, Ukraine
Tel.: (8032)239 45 31 (off.).