Cyril REGAN

May 26, 2021



## Table of Contents

### Reinforcement Learning

Markov Decision Process (MDP) Model Free - Model Based Dyna Model

### Successor Representation

me based representation SR matrix construction SR relation with transition function SR policy dependence

### Replay

ΕV

Bellman backur

Gain - Need
Algorithm limitations
Bidirectionnal MB
Prioritized Sweeping
Trajectory Sampling
Bidirectional Algorith

#### LSFM

State of the Art Literature Latent space Latent space an

Latent space and Successor Fea-

tures Learning

Conclusion







## Markov Decision Process (MDP)

15

RL purpose is to solve Markov Decision Process (MDP)  $\langle S, A, p, r, \gamma \rangle$ 

 $\mathcal{S}$ : State Space :  $s \in \mathcal{S}$ 

A: Action Space :  $a \in A$ 

Transition function p is Markovian : p(s'|s, a) = $\mathbb{P}[s_{t+1} = s' | s_t = s, a_t = a]$ 

: Reward function  $r: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [\mathbb{R}, \mathbb{R}]$ 

discount factor that captures a preference for prox-

imal rewards







But what is solve MDP ? Find the *optimal policy*  $\pi^*$  to maximize the rewards. The *policy*  $\pi$  is a function that maps the state of an agent to an action :



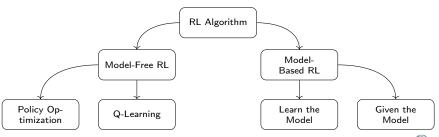
$$\pi(a|s) = \mathbb{P}[a_t = a|s_t = s]$$





## 2 main type of Models to solve a MDP:

- Model-based : An internal model (p, r) is known or learned
- Model-free : No internal model







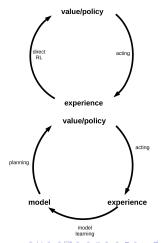
### Model Free - Model Based

#### Model Free:

- Need many on-line steps,
- Inflexible (transition / reward change),
- + No planning phase,
- + Computationally cheap.

## Model based:

- Need to learn model (can be given),
- + few on-line steps (off-line learning),
- + flexible (transition / reward change),
- planning can be complex (trajectory optimization),
- computationally expensive.

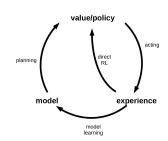


# Dyna Model

### Dyna Model:

- Model Free learning on-line.
- Model Based Planning off-line

Simulated experiences similar to hippocampal replay



Dyna Model Based planning can be constructed with the :

## Successor Representation (SR)

(instead to the transition function p).



## Table of Contents

#### Reinforcement Learning

Markov Decision Process (MDP Model Free - Model Based Dyna Model

### Successor Representation

me based representation
SR matrix construction
SR relation with transition function
SR policy dependence
TD-Learning

### Replay

EVI

Bellman backup

Gain - Need
Algorithm limitations
Bidirectionnal MB
Prioritized Sweeping
Trajectory Sampling
Bidirectional Algorithm

#### LSFM

State of the Art Literature Latent space Latent space and

Latent space and Successor Fea-

tures Learning

Conclusion







# Successor Representation: a time based representation

SR matrix construction

(Dayan 1993; Gershman 2018)

SR is a "compressed time based" representation of the actual state and all the future state occupancy (discounting by the  $\gamma$  factor):

$$\mathbf{M}^{\pi}(s', s, a) = \mathbb{E}_{\pi}\left[\sum \gamma^{t} \mathbb{I}\{s_{t} = s'\} | s_{0} = s, a_{0} = a\right]$$
 (1)

Where  $\mathbb{I}\{...\}=1$  if its argument is true (0 otherwise). In tabular state, SR matrix  $\mathbf{M}^{\pi}$  relation with transition matrix  $(\mathbf{\mathcal{P}}^{\pi})_{ss'}^{a} = p(s'|s,a)$  is :

$$\mathbf{M}^{\pi}(s', s, a) = \mathbb{I}\{s = s'\} + \gamma(\mathbf{P}^{\pi})_{ss'}^{a} + \gamma^{2}(\mathbf{P}^{\pi 2})_{ss'}^{a} + \gamma^{3}(\mathbf{P}^{\pi 3})_{ss'}^{a} + \dots (2)$$

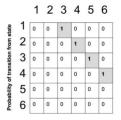
In a matrix formulation :  $\mathbf{M}^{\pi} = (1 - \gamma \mathbf{\mathcal{P}}^{\pi})^{-1}$ 

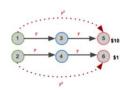


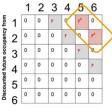


### SR relation with transition function

SR Matrix  $\mathbf{M}^{\pi}$ : a time based representation:







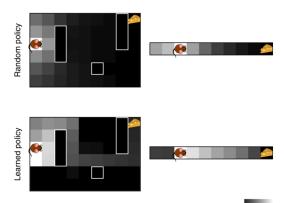
Transition matrix  $(\mathcal{P}^{\pi})$  Mdp (Momennejad and Howard 2018)







## SR policy dependence



Successor Representation of a state (row of M)

Successor Representation policy dependence (Mattar-Daw 2018)



Reinforcement Learning

SR Matrix  $\mathbf{M}^{\pi}$  can be learn as TD-Learning (Gershman 2018):

$$\boldsymbol{M}^{\pi}(s,:) \leftarrow \boldsymbol{M}^{\pi}(s,:) + \alpha(1_s + \gamma \boldsymbol{M}^{\pi}(s',:) - \boldsymbol{M}^{\pi}(s,:))$$
 (4)

Nice property: Value functions (V, Q) are **linear** with SR (M):

$$V^{\pi}(s) = \mathbb{E}_{\pi} \left[ \sum_{k=0}^{+\infty} \gamma^k r_{t+k} | s_t = s \right] = \sum_{s'} \mathbf{M}^{\pi}(s, s') \mathbf{r}(s') \tag{5}$$

$$Q^{\pi}(s,a) = \mathbb{E}_{\pi}\left[\sum_{k=0}^{+\infty} \gamma^k r_{t+k} | s_t = s, a_t = a\right] = \sum_{s'} \mathbf{M}^{\pi}(s,s',a) \mathbf{r}(s',a)$$

## Table of Contents

### Reinforcement Learning

Markov Decision Process (MDP Model Free - Model Based Dyna Model

### Successor Representation

Ime based representation
SR matrix construction
SR relation with transition function
SR policy dependence

### Replay

EV

Bellman backup

Gain - Need
Algorithm limitations
Bidirectionnal MB
Prioritized Sweeping
Trajectory Sampling
Bidirectional Algorithm

#### LSFM

State of the Art
Literature
Latent space
Latent space and S

Latent space and Successor Fea-

tures Learning

Conclusion





- Model-Free: experience storage + off-line learning (Lin 1992)
  - Random experience replay (Mnih 2013; Mnih 2015)
  - prioritized experience replay (Schaul 2016)
- Model-Based / SR-Based : simulated experience replay
- $2 \neq \text{prioritized simulated replay method}$ :
  - The Expected Value of a Bellman backup (EVB) (Mattar-Daw 2018)
  - The Bidirectionnal Model-Based approach (Khamassi-Girard 2020).



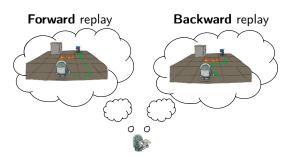
## 2 computational families of simulated replay:

EVB Need Gain (Mattar-Daw 2018):

Bidirectionnal MB (Khamassi-Girard 2020):

Trajectory sampling

Prioritized sweeping



## Expected Value of a Bellman backup (EVB)

## The Bellman backup

EVB (Mattar-Daw 2018) algorithm is based on :

- Dyna architecture,
- and prioritized Bellman backup.

A Bellman backup is the TD-learning iteration of Q:

$$Q(s_k, a_k) \leftarrow Q(s_k, a_k) + \alpha \left[ r_k + \gamma \max_{a} Q(s'_k, a) - Q(s_k, a_k) \right]$$
 (7)

A Bellman backups occur on-line or off-line with simulated replay.





# EVB Policy :

Softmax Q function 
$$\pi(a|s) = \frac{e^{\beta Q(s,a)}}{\sum_{i} e^{\beta Q(s,i)}}$$
:



Pure **exploration** 

 $\beta \rightarrow 0^+$ 



= Probability for all actions

Probability = 1 for the action with max Q value (= 0 for the others)



0000000

# Expected Value of a Bellman backup (EVB):

EVB is the criterion to  $\underline{\text{order}}$  the replay experiences for the offline update of Q (Bellman Backup).

$$EVB(s_k, a_k) = Gain(s_k, a_k) \times Need(s_k)$$
 (8)

- **Gain**: expected return when reaching  $(s_k, a_k)$ ,
- **Need**: number of times the agent expects to visit  $s_k$  in the future.





Gain

$$Q_{\pi}(\mathsf{s}_k, \mathsf{a}_k) = \mathbb{E}_{\pi}\left[\sum_{k}^{\infty} \gamma^i r_{t+i}
ight]$$

**Gain** on  $(s_k, a_k)$  quantifies the **improvement** in the expected return  $(\pi_{new}(a|s_k) - \pi_{old}(a|s_k))$  $\boldsymbol{Q}^{new}(s_k,a_k) \leftarrow \boldsymbol{Q}^{old}(s_k,a_k)$ 

as a result of a **policy change** after a Bellman Backup on  $(s_k, a_k)$ :

$$Gain(s_k, a_k) = \sum_{a} \overbrace{Q_{\pi_{new}}(s_k, a)}^{\text{expected return}} \left( \pi_{new}(a|s_k) - \pi_{old}(a|s_k) \right)$$
(9)

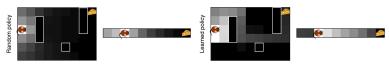
where  $\pi_{new}(a|s_k) = softmax(\mathbf{Q}^{new}(a,s_k))$  (resp.  $\pi_{old}$ )



### Need

**Need** on  $s_k$  quantifies the discounted number of times the agent expects to visit states in the future :

$$Need(s_k) = \sum_{i}^{\infty} \gamma^i.1_{S_{t+i,s_k}} = \sum_{i \in \mathcal{S}} M[s_k, s_i] = \sum_{row \ s_k} M$$
 (10)



 $\blacksquare$  : SR of a state (row of M)



### EVB algorithm and limitations :

Offline learning : Compute fixed number of planning steps on on the start and the end of each episode (not flexible).

- planning step :
  - \* EVB is computed on all possible experiences (very high computational cost).
  - \* Q update (Bellman backup) on experience with the higher EVB

 $\underline{\text{Online learning}} = Q \text{ learning}: \text{Update (Bellman backup) on current experience}$ 



## EVB Algorithm:

$$(s_k, a_k, r_k, s'_k)_{k \in \mathcal{S} \times \mathcal{A}}$$

Loop on all states-actions and store all experiences of the environment

Construct Transition ( $\mathcal{P}$ ) and SR ( $M = (1 - \gamma \mathcal{P})^{-1}$ ) matrices

loop on  $n_{max}$  episodes while  $n <= n_{max}$  do

 $p_{max}$  offline learning steps at the start of episode while  $p <= p_{max}$  do

Compute EVB on all states-actions  $(s_k, a_k)$ 

Belman Backup  $(s_k, a_k, r_r, s'_k)$ 

UpdateQ with the experience with higher EVB

#### end

#### online learning

$$a_t \leftarrow argmax_a softmax_\beta(Q(s_t, a))$$

Take action  $a_t$  receive  $(s_{t+1}, r_{t+1}, done)$ 

Update Q (Belman Backup),  ${\cal P}$  and  ${\it M}$  on  $(s_t,a_t,s_{t+1},r_{t+1})$ 

 $p_{max}$  offline learning steps  $\underline{again}$  at the  $\underline{end}$  of episode while  $p <= p_{max}$  do

end

end



# Bidirectionnal Model-Based (MB)

Bidirectionnal MB : prioritized replay with a heuristic-based approach (Khamassi-Girard 2020) .

- Backward replay : Prioritized sweeping (PS)
  - \* Start : replay starting from states whose Q-value changed recently,
  - \* and then to their predecessors
  - \* and the predecessors to their predecessors
  - ... and so on.
- Forward replay : Trajectory sampling (TS)
  - $\ast$  Start : generate continuous trajectories from the current position
  - ... until forward and backward search reached a connection state.





PS is backward replay. PS consists of weighing experiences based on how surprising they are, measured by the absolute value of the MB TD error signal:

$$\delta = |\boldsymbol{Q}^{n+1}(s,a) - \boldsymbol{Q}^{n}(s,a)|,$$
  
$$\boldsymbol{Q}^{n+1}(s,a) \leftarrow r(s,a) + \gamma \sum_{s'} \mathcal{P}_{ss'}^{a} \max_{a'} \boldsymbol{Q}^{n}(s',a')$$
 (11)



- equal to  $\delta$ ,
- attenuated by  $\gamma$ ,
- and attenuated by the probability to effectively reach state s from the predecessor s' under consideration : T(s', a', s)

$$p \leftarrow \delta \times \gamma \times T(s', a', s) \tag{12}$$



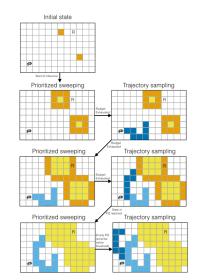
- TS drives forward replay from the current location.
- TS generates trajectories(with transition state function).
- TS stops when a connection state with backward search (Prioritized Sweeping) is reached.

## Bidirectional Algorithm



### Offline alternation between

- Prioritized Sweeping,
- Trajectory Sampling.





loop until 
$$n_{a_{max}}$$
 actions while  $n_a = n_{a_{max}}$  do 
$$a_t \leftarrow softmax_\beta(Q(s_t, a_t)) \mid \text{Take action } a_t \text{ receive } s_{t+1}, r_t,$$
 Update  $\mathcal{P}_{s_t s_{t+1}}^{a_t} \mid R(s_t, a_t) \leftarrow r_t \mid \text{Store } Q_{old} \leftarrow Q(s_t, a_t)$  
$$Q(s_t, a_t) \leftarrow R(s_t, a_t) + \gamma \sum_{s'} \mathcal{P}_{s_t s_{t+1}}^{a_t} \max_k Q(s', k)$$
 Compute the level of surprise:  $\delta = |Q(s_t, a_t) - Q_{old}|$  reach predecessors  $(s, a)$  of  $s_t$  for each  $(s, a)$  so that  $\mathcal{P}_{ss_t}^a \neq 0$  do Compute predecessors priority:  $p \leftarrow \delta \times \gamma \times \mathcal{P}_{ss_t}^a$  if  $(s, a) \notin PQ$ ueue then: Put  $(s, a)$  in  $PQ$ ueue with priority  $p$  else: Update priority of  $(s, a)$  in  $PQ$ ueue with  $p$  end 
$$s_t \leftarrow s_{t+1} \text{ and } n_a = n_a + 1$$
 Offline MB  $Q$ -update  $Q \leftarrow OfflineBidirectional( $Q, PQ$ ueue,  $s_t$ ) end$ 

Algorithm 1: Online bidirectional algorithm



$$nbLoops \leftarrow 0$$
 repeat

```
Prioritized Sweeping (PS) while PQueue[0] > thresold | max budget do
```

 $\mathsf{Select}\ \underline{\mathbf{highest}\ \mathbf{prioritized}}\ \mathbf{tuple}\ \mathbf{:}\quad (s,a) \leftarrow PQueue[0]$ 

MB **Q Update** on Q(s, a)

Update (s, a) <u>predecessors</u> priority : for each (s', a') ... do ...

### end

TS start from the **current location** in the environment :  $s \leftarrow s_t$ 

**Trajectory Sampling (TS)** while  $s \notin PQueue \mid max \ budget \ do$ 

$$a \leftarrow softmax_{\beta}(Q(s)) \mid s' \leftarrow the \ \textit{best probability of } T(s, a, s')$$

$$Q_{old} \leftarrow Q(s, a) \mid Q(s, a) \leftarrow R(s, a) + \gamma \sum_{s'} \mathcal{P}_{ss'}^{a} \max_{k} Q(s', k)$$

Sum of **TD error** : 
$$Sum\delta \leftarrow Sum\delta + |Q(s,a) - Q_{old}|$$

Jump to the **next state** :  $s \leftarrow s'$ 

end

until TD error convergence or budget max; return Q



## Table of Contents

### Reinforcement Learning

Markov Decision Process (MDP) Model Free - Model Based Dyna Model

### Successor Representation

me based representation
SR matrix construction
SR relation with transition function
SR policy dependence
TD-Learning

### Replay

EVB

Bellman backur

Gain - Need
Algorithm limitations
Bidirectionnal MB
Prioritized Sweeping
Trajectory Sampling
Bidirectional Algorith

#### **LSFM**

State of the Art Literature Latent space Latent space and

Latent space and Successor Fea-

tures Learning

Conclusion



### State of the Art

SR can represent small discrete Tabular State Space => Unsuitable for big states space like video-games environment.



- => Need approximate functions to reduce the original state space
- => The resulting reduced **latent** representation can lead to **faster learning** because information can be **re-used** across different inputs





$$\mathbf{M}^{\pi}(s', s, a) = \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^{t} \mathbb{I}\{s_{t} = s'\} | s_{0} = s, a_{0} = a \right]$$
 (13)

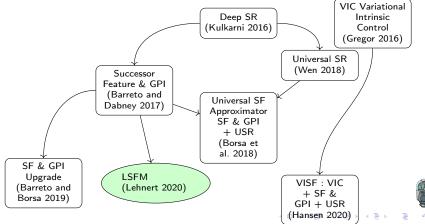
$$\psi^{\pi}(s,a) = \mathbb{E}_{\pi} \left[ \sum_{t=1}^{\infty} \gamma^{t-1} \phi_{st} \middle| s_1 = s, a_1 = a \right], \tag{14}$$

 $\pmb{\psi}^{\pi}$  : frequency of different latent states vectors encountered following  $\pi.$ 



### Literature

Some papers on SR Approximation functions in Neural Network.



## Latent space

Latent state space  $S_\phi$  is constructed using a state representation function :

$$\phi: S \to S_{\phi} \tag{15}$$

•00

 $S_{\phi}$  is  $rac{ ext{value-predictive}}{ ext{reward-predictive}}$  if it retains enough information to support

$$V^{\pi}$$
 or  $Q^{\pi}$   $\left(\mathbb{E}_{\pi}\left[\sum_{k}\gamma^{k}r_{k}\right]\right)$ 

- accurate value or Q-value predictions
- accurate future expected reward predictions

$$\mathbb{E}[r_1,r_2,\ldots|s,a_1,a_2,\ldots]$$



- Value-predictive state representation is policy dependent.
- Reward-predictive state representation is not.

For optimal policy  $\pi^*$ , value-predictive state representation becomes a

Reward-maximizing state representation



(s, a, r, s') experiences on  $s_0 = s_{21}, \rightarrow, \uparrow, \rightarrow$  with Reward Predictive and Maximizing state representations (Lucas Lehnert n.d.)

Actions:  $\rightarrow \uparrow \leftarrow \downarrow$ 

Compressed Task:

Original Task:

(A) Column World Task

 $\phi_1\phi_2\phi_3$ 



(B) Reward-Predictive State Representation





(C) Reward-Maximizing State Representation

#### Original State

### Reward-Predictive

Reward-Maximizing





Reward-maximizing state representations of (Barreto and Dabney 2017) are directly the SFs ( $\psi^{\pi}$ ):

$$Q^{\pi}(s,a) = \psi^{\pi}(s,a).\mathbf{w} \tag{16}$$

Latent space  $\phi$  is one-step rewards :  $\mathbf{r}(s, a, s') = \phi(s, a, s') \cdot \mathbf{w}$ .

**Reward-predictive** state representations  $(\phi_s)$  of **Linear Successor Features Model** (LSFM) are tied with SFs ( $\psi^{\pi}$ ) via linear matrices  $\{\boldsymbol{F}_a\}_{a\in\mathcal{A}}$  (Lehnert 2020) :

$$\phi_s^{\top} \mathbf{F}_a \approx \psi^{\pi}(s, a),$$
 (17)



$$\mathbf{F}_{a} = \mathbf{I} + \gamma \mathbf{M}_{a} \overline{\mathbf{F}}, \text{ where } \overline{\mathbf{F}} = \frac{1}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} \mathbf{F}_{a}$$
 (18)

000

A LAM  $\{M_a\}_{a\in\mathcal{A}}$  model the empirical latent transition probabilities :

$$\phi_s^{\top} \mathbf{M}_a = \mathbb{E}[\phi_s'|s,a], \quad \text{and} \quad \phi_s^{\top} \mathbf{w}_a = \mathbb{E}[r(s,a,s')|s,a]$$
 (19)

 $\phi$  is (multi) **Reward-predictive** :

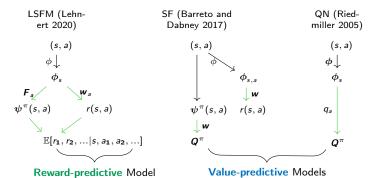
$$\mathbb{E}[r_t|s,a_1,...,a_t] \approx \phi_s^{\top} \boldsymbol{M}_{a1}...\boldsymbol{M}_{at-1}.\boldsymbol{w}_a$$
 (20)



References

000

#### Value and Reward state representation models :







## Learning

#### LSFM Learning

$$\mathcal{L}_{LSFM} = \underbrace{\sum_{i=1}^{D} \left(\phi_{s_{i}}^{\top} \mathbf{w}_{a_{i}} - r_{i}\right)^{2}}_{=\mathcal{L}_{r}} + \alpha_{\psi} \underbrace{\sum_{i=1}^{D} \left\|\phi_{s_{i}}^{\top} \mathbf{F}_{a} - \vec{\mathbf{y}}_{s_{i}, a_{i}, r_{i}, s_{i}'}\right\|_{2}^{2}}_{=\mathcal{L}_{\psi}} + \underbrace{\alpha_{N} \sum_{i=1}^{D} \left(\left\|\phi_{s_{i}}\right\|_{2}^{2} - 1\right)^{2}}_{=\mathcal{L}_{N}}$$
(21)

and  $\alpha_N$ ,  $\alpha_{\psi} > 0$  are hyper-parameters.

$$\vec{\mathbf{y}}_{s,a,r,s'} = \boldsymbol{\phi}_s^\top + \gamma \boldsymbol{\phi}_{s'}^\top \overline{\mathbf{F}}$$
 (22)

•000

 $\overline{\emph{\textbf{F}}}$  the LSFM that predicts the SF for a policy  ${f uniformly}$  at  ${f random}$  :

$$\overline{\mathbf{F}} = \frac{1}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} \mathbf{F}_a \tag{23}$$



A theorem (Lehnert 2020) bounds  $\underline{\text{the prediction error}}$  of finding a linear approximation of the Q-function :

$$Q^{\pi}(s,a)pprox \phi_s^{ op}oldsymbol{q}_a$$

using a state representation  $\phi$  and a real valued vector  $\boldsymbol{q}_a$ .



## Learned Representation (Lehnert 2020)

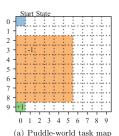
Puddle Word /

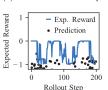
Agglomerating clustering for LSFM learned representation

Expected reward and predictions on random action sequence

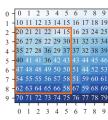
Random representation /

LSFM representation



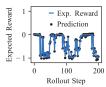


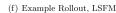




0000

(c) LSFM generalization map





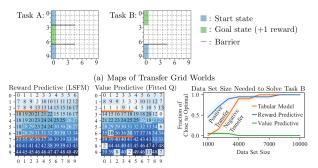






000

## Learned transfer (Lehnert 2020)



(c) State partitions obtained through learning on a transition data set D.

LSFM representation allow **faster transfer learning** than tabular model-based agent and Value Predictive representation.



value iteration construction on optimal policy

## Table of Contents

#### Reinforcement Learning

Markov Decision Process (MDP) Model Free - Model Based Dyna Model

#### Successor Representation

me based representation
SR matrix construction
SR relation with transition function
SR policy dependence

## Replay

EVE

Bellman backur

Gain - Need
Algorithm limitations
Bidirectionnal MB
Prioritized Sweeping
Trajectory Sampling
Bidirectional Algorith

#### **LSFM**

State of the Art
Literature
Latent space
Latent space and

Latent space and Successor Fea-

tures Learning

Conclusion



## Conclusion & future work Conclusion

Reinforcement Learning

#### LSFM is **Reward-predictive**

LSFM state representations generalize across variations in transition and reward functions (in contrast to (Barreto and Dabney 2017), (Kulkarni 2016), ... in reward functions only).

LSFMs must be trained on pure exploration policy (reward agnostic).

The LSFM state representation can be used to predict the value-function of any policy, including the optimal policy.



#### Futur work

\*\*\*

# <u>Goal</u>: Construction of a RL model using <u>offline generated</u> Replay with LSFM Successor Features

for Video games (and Neuroscience ?) applications.

\*

## Replay State of the Art:

(Mattar-Daw 2018) and (Khamassi-Girard 2020) proposed mathematical and heuristic models to generate Replay in DYNA algorithm on <u>original</u> tabular state space.



#### Research directions:

- => Find new or adapt existing **Replay** models on Successor Features
- => Find better (than random) **exploration policy** for LSFM learning : *Options* of (Machado 2018) ?

Any ideas and helps is welcome

Thanks for your attention!



# References I



- Barreto, Andre and Will Dabney (2017). Successor Features for Transfer in Reinforcement Learning. 31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA.
  - Borsa, Diana et al. (2018). "universal successor feature approximator." In:
- Dayan, Peter (1993). Improving Generalization for Temporal Difference Learning: The Successor Representation. Computational Neurobiology Laboratory, The Salk Institute, P.O. Box 85800, Sun Diego, CA 92186-5800 USA.
  - Gershman, Samuel J. (2018). *The Successor Representation: Its Computational Logic and Neural Substrates*. The Journal of Neuroscience. ISBN: 9781417642595.



## References II



- Hansen, Barreto (2020). "Fast Task Inference with Variational Intrinsic Successor Features." In:
- Khamassi-Girard (2020). Modeling awake hippocampal reactivations with model-based bidirectional search. Biological Cybernetics (Modeling), Springer Verlag, 2020, 10.1007/s00422-020- 00817-x. hal-02504897.
- Kulkarni, Gershman (2016). Deep Successor Reinforcement Learning. ArXiv vol abs/1606.02396.
- Lehnert, Lucas (2020). Successor Features Combine Elements of Model-Free and Model-based Reinforcement Learning.
  - Lin (1992). Self-improving reactive agents based on reinforcement learning, planning and teaching. Machine learning, 8(3):293–321, 1992.



## References III





- Mattar-Daw (2018). Prioritized memory access explains planning and hippocampal replay. Nature Neuroscience.
  - Mnih (2013). Playing atari with deep reinforcement learning. ICLR 2016arXiv:13 2013. Mnih.
    - (2015). Human-level control through deep reinforcement learning. Nature, 518(7540):529-533, 2015.
  - Momennejad, Ida and Marc W. Howard (Oct. 22, 2018). Predicting the Future with Multi-scale Successor Representations. preprint. Neuroscience. DOI: 10.1101/449470. URL: http://biorxiv.org/ lookup/doi/10.1101/449470 (visited on 04/21/2021).

## References IV



Riedmiller, Martin (2005). "Neural Fitted Q Iteration – First Experiences with a Data Efficient Neural Reinforcement Learning Method." en. In: Machine Learning: ECML 2005. Ed. by David Hutchison et al. Vol. 3720. Series Title: Lecture Notes in Computer Science. Berlin, Heidelberg: Springer Berlin Heidelberg, pp. 317–328. ISBN: 978-3-540-29243-2 978-3-540-31692-3. DOI: 10.1007/11564096\_32. URL: http://link.springer.com/10.1007/11564096\_32 (visited on 04/30/2021).



Schaul (2016). Prioritized experience replay. ICLR.



Wen, Chen Ma & Junfeng (2018). "universal successor representations for transfer reinforcement learning." In: Workshop track - ICLR 2018.

