



# DOES ZERO-SHOT REINFORCEMENT LEARNING EXIST?

Ahmed Touati, Jérémy Rapin & Yann Ollivier Meta AI Research, Paris, {atouati, jrapin, yol}@meta.com

**ICLR 2023** 

# Definition



What is Zero-Shot Reinforcement Learning?

**Zero-shot RL: problem statement.** The goal of zero-shot RL is to compute a compact representation  $\mathcal{E}$  of the environment by observing samples of reward-free transitions  $(s_t, a_t, s_{t+1})$  in this environment. Once a reward function is specified later, the agent must use  $\mathcal{E}$  to immediately produce a good policy, via only elementary computations without any further planning or learning. Ideally, for any downstream task, the performance of the returned policy should be close to the performance of a supervised RL baseline trained on the same dataset labeled with the rewards for that task.

Loose Limitations: How to define Zero-Shot RL with planning assisted?



# Motivation

ParNeC 模式识别与神经计算研究组 PAttern Recognition and NEural Computing

Why Zero-Shot Reinforcement Learning?







#### Offline goal-conditioned Reinforcement Learning with ICM-like exploration?

Advantages: Easy to Implement.

Disadvantages: goal-conditioned policy learned depends on the behavior policy.

#### Diffusion Model for trajectories with rewards conditioned?

Advantages: No need for data preprocess and subtle design. Disadvantages: Low interpretability and solidity.

### Goal-conditioned or Skill-based Planning with landmarks?

Advantages: Long-range planning and High interpretability. Disadvantages: Consists of too many components and cost expensive.

### Q values and V values decoupled into successor representation and reward representation?

Advantages: Easy to implement and Solid theoretical analysis Disadvantages: Reward Space is infinite in continuous state space. Perspective: View different tasks as different reward functions motivated.

# **Decouple-Successor Features**

a0



$$M \equiv (S, A, p, r, \gamma)$$

$$izide (S, A, p, r, \gamma)$$

$$Q^{\pi}(s, a) \equiv \mathbb{E}^{\pi} [G_t | S_t = s, A_t = a] \quad \text{where} \quad G_t = \sum_{i=0}^{\infty} \gamma^i R_{t+i}$$

$$\pi'(s) \in \underset{a}{\operatorname{argmax}} Q^{\pi}(s, a);$$

$$Q^{\pi}(s, a) = \mathbb{E}^{\pi} [\sum_{i=t}^{\infty} \gamma^{i-t} \phi_{i+1} | S_t = s, A_t = a]^{\top} \mathbf{w} \equiv \psi^{\pi}(s, a)^{\top} \mathbf{w}, \text{ where } r(s, a, s') = \phi(s, a, s')^{\top} \mathbf{w},$$

$$griede (S, A, p, r, \gamma)$$

$$Q^{\pi}(s, a) \equiv \mathbb{E}^{\pi} [G_t | S_t = s, A_t = a] \quad \text{where} \quad G_t = \sum_{i=0}^{\infty} \gamma^i R_{t+i}$$

$$\pi'(s) \in \underset{a}{\operatorname{argmax}} Q^{\pi}(s, a);$$

$$Q^{\pi}(s, a) = \mathbb{E}^{\pi} [\sum_{i=t}^{\infty} \gamma^{i-t} \phi_{i+1} | S_t = s, A_t = a]^{\top} \mathbf{w} \equiv \psi^{\pi}(s, a)^{\top} \mathbf{w}, \text{ where } r(s, a, s') = \phi(s, a, s')^{\top} \mathbf{w},$$

$$griede (S, A, p, r, \gamma)$$

$$Q^{\pi}(s, a) \equiv \mathbb{E}^{\pi} [\sum_{i=t}^{\infty} \gamma^{i-t} \phi_{i+1} | S_t = s, A_t = a]^{\top} \mathbf{w} \equiv \psi^{\pi}(s, a)^{\top} \mathbf{w}, \text{ where } r(s, a, s') = \phi(s, a, s')^{\top} \mathbf{w},$$

$$griede (S, A, p, r, \gamma)$$

$$griede (S, A, p, r, \gamma)$$

$$Q^{\pi}(s, a) \equiv \mathbb{E}^{\pi} [\sum_{i=t}^{\infty} \gamma^{i-t} \phi_{i+1} | S_t = s, A_t = a]^{\top} \mathbf{w} \equiv \psi^{\pi}(s, a)^{\top} \mathbf{w}, \text{ where } r(s, a, s') = \phi(s, a, s')^{\top} \mathbf{w},$$

$$griede (S, A, p, r, \gamma)$$

$$griede (S, A, p, r, \gamma)$$

$$Q^{\pi}(s, a) \equiv \mathbb{E}^{\pi} [\sum_{i=t}^{\infty} \gamma^{i-t} \phi_{i+1} | S_t = s, A_t = a]^{\top} \mathbf{w} \equiv \psi^{\pi}(s, a)^{\top} \mathbf{w}, \text{ where } r(s, a, s') = \phi(s, a, s')^{\top} \mathbf{w},$$

$$griede (S, A, p, r, \gamma)$$

$$griede (S, A, p, r, \gamma)$$

$$Q^{\pi}(s, a) \equiv \mathbb{E}^{\pi} [\sum_{i=t}^{\infty} \gamma^{i-t} \phi_{i+1} | S_t = s, A_t = a]^{\top} \mathbf{w} \equiv \psi^{\pi}(s, a)^{\top} \mathbf{w}, \text{ where } r(s, a, s') = \phi(s, a, s')^{\top} \mathbf{w},$$

$$griede (S, A, p, r, \gamma)$$

$$griede (S, A, p, r, \gamma)$$

$$Q^{\pi}(s, a) \equiv \mathbb{E}^{\pi} [\sum_{i=t}^{\infty} \gamma^{i-t} \phi_{i+1} | S_t = s, A_t = a]^{\top} \mathbf{w} \equiv \psi^{\pi}(s, a)^{\top} \mathbf{w}, \text{ where } r(s, a, s') = \phi(s, a, s')^{\top} \mathbf{w}, \text{ where } r(s, a, s') = \phi(s, a, s')^{\top} \mathbf{w}, \text{ where } r(s, a, s') = \phi(s, a, s')^{\top} \mathbf{w}, \text{ where } r(s, a, s') = \phi(s, a, s')^{\top} \mathbf{w}, \text{ where } r(s, a, s') = \phi(s, a, s')^{\top} \mathbf{w}, \text{ where } r(s, a, s') = \phi(s, a, s')^{\top} \mathbf{w}, \text{ where } r(s, a, s') = \phi(s, a, s')^{\top} \mathbf{w}, \text{ where } r(s, a, s') = \phi(s, a, s')^{\top} \mathbf{w}, \text{ where } r(s, a, s') = \phi(s, a, s')^{\top} \mathbf{w}, \text{ where } r(s, a, s') = \phi(s, a,$$













ParN,C 模式识别与神经计算研究组 PAttern Recognition and NEural Computing







# How to learn Successor Features



# FAST TASK INFERENCE WITH VARIATIONAL INTRINSIC SUCCESSOR FEATURES

$$Q^{\pi}(s, a) = \boldsymbol{\psi}(s, a, e(\pi))^{\top} \mathbf{w}.$$

$$\begin{aligned} \mathcal{F}(\theta) &= I(z; f(\tau_{\pi_{\theta}})) = H(z) - H(z|f(\tau_{\pi_{\theta}})). \\ \mathcal{F}(\theta) &= -H(z|f(\tau_{\pi_{\theta}})). \\ \mathcal{F}(\theta) &= \sum_{s,z} p(s,z) \log p(z|s) = \mathbb{E}_{\pi,z} [\log p(z|s)]. \\ -H(z|s) &= \sum_{s,z} p(s,z) \log p(z|s) \\ &= \sum_{s,z} p(s,z) \log p(z|s) + \sum_{s,z} p(s,z) \log q(z|s) - \sum_{s,z} p(s,z) \log q(z|s) \\ &= \sum_{s,z} p(s,z) \log q(z|s) + \sum_{s,z} p(s,z) \log p(z|s) - \sum_{s,z} p(s,z) \log q(z|s) \\ &= \sum_{s,z} p(s,z) \log q(z|s) + \sum_{s} p(s) KL(p(\cdot|s)) ||q(\cdot|s)) \\ &\geq \sum_{s,z} p(s,z) \log q(z|s) \\ &= \mathbb{E}_{\pi,z} [\log q(z|s)] \end{aligned}$$
(12)

 $\pi(s) = \underset{a}{\operatorname{argmax}} \max_{i} \psi(s, a, e(\pi_i))^{\top} \mathbf{w} = \underset{a}{\operatorname{argmax}} \max_{i} Q^{\pi_i}(s, a).$ 

# How to learn Successor Features





Figure 1: VISR model diagram. In practice  $\mathbf{w}_t$  is also fed into  $\boldsymbol{\psi}$  as an input, which also allows for GPI to be used (see Algorithm 1 in Appendix). For the random feature baseline, the discriminator q is frozen after initialization, but the same objective is used to train  $\boldsymbol{\psi}$ .



Aimed at maximizing the reward collected policy under reward z (also w)

 $r(s, a, s') = \boldsymbol{\phi}(s, a, s')^{\top} \mathbf{w},$ 

What is predefined, View reward as a fixed point we can denote that.

$$\log q(\mathbf{w}|s) = \boldsymbol{\phi}(s)^T \mathbf{w}.$$

# How to learn Successor Features



Algorithm 1: Training VISR Randomly Initialize  $\phi$  network // L2 normalized output layer //  $dim(output) = #A \times dim(W)$ Randomly Initialize  $\psi$  network for  $e := 1, \infty$  do sample w from L2 normalized  $\mathcal{N}(0, I(dim(W)))$ // uniform ball  $Q(\cdot, a|w) \leftarrow \psi(\cdot, a, w)^\top w, \forall a \in A$ for t := 1, T do Receive observation  $s_t$  from environment if using GPI then  $Q^{\pi_0}(\cdot, a) \leftarrow Q(\cdot, a|w) \forall a \in A$ for i := 1, 10 do sample  $w_i$  from VMF( $\mu = w, \kappa = 5$ )  $Q^{\pi_i}(\cdot, a) \leftarrow \psi(\cdot, a, w_i)^\top w, \forall a \in A$ end  $a_t \leftarrow \epsilon$ -greedy policy based on  $\max_i Q^{\pi_i}(s_t, \cdot)$ else  $a_t \leftarrow \epsilon$ -greedy policy based on  $Q(s_t, \cdot | w)$ end Take action  $a_t$ , receive observation  $s_{t+1}$  from environment  $a' = \operatorname{argmax}_{a} \psi(s_{t+1}, a, w)^{\top} w$  $y = \phi(s_t) + \gamma \psi(s_{t+1}, a', w)$  $loss_{\psi} = \sum_{i} (\psi_i(s_t, a_t, w) - y_i)^2$  $loss_{\phi} = -\phi(s_t)^{\top} w$ // minimize Von-Mises NLL Gradient descent step on  $\psi$  and  $\phi$ // minibatch in practice end end

# **Decouple-Successor States**



## What M really is and the generalization



$$M^{\pi}(s_0, a_0, X) := \sum_{t \ge 0} \gamma^t \Pr\left((s_t, a_t) \in X \mid s_0, a_0, \pi\right)$$

Assume a goal-conditioned env where the goal is  $s_t$ ,  $a_t$  from initial state  $s_0$ ,  $a_0$ . M is equal to Q( $s_t$ ,  $a_t$ ). So M is actually the discount distance between arbitrary state-action pairs under policy  $\pi$ .



If we want s1,a1 from initial s0, We can obtain from the M: Q(s0,a0)<Q(s0,a1), So we need to choose a1.

IIINotably, the M is conditioned by the random policy, but can generalize to the goal-conditioned tasks.





Like the calculation of Q value, we can use dynamic process to calculate M.

Backward: 
$$M = \mathcal{I} + \gamma MP$$
 Forward:  $M = \mathcal{I} + \gamma PM$ 

Define a target update of M via the Bellman equation,  $M^{\text{tar}} := \text{Id} + \gamma P M_{\theta_t}$ . Define the loss between M and  $M^{\text{tar}}$  via  $J(\theta) := \frac{1}{2} \|M_{\theta} - M^{\text{tar}}\|_{\rho}^2$  using the norm (1). Then the gradient step on  $\theta$  to reduce this loss is

 $\widehat{\delta\theta_{\rm BTD}} = \partial_{\theta} \tilde{m}_{\theta}(s,s) + \tilde{m}_{\theta}(s_1,s) \left( \gamma \, \partial_{\theta} \tilde{m}_{\theta}(s_1,s') - \partial_{\theta} \tilde{m}_{\theta}(s_1,s) \right)$ 

$$- \partial_{\theta} J(\theta)_{|\theta=\theta_{t}} = \mathbb{E}_{s \sim \rho, s' \sim P(s, \mathrm{d}s'), s_{2} \sim \rho} \left[ \partial_{\theta} \tilde{m}_{\theta_{t}}(s, s) + \partial_{\theta} \tilde{m}_{\theta_{t}}(s, s_{2}) \left( \gamma \tilde{m}_{\theta_{t}}(s', s_{2}) - \tilde{m}_{\theta_{t}}(s, s_{2}) \right) \right].$$
(21)

### **Matrix Factorization**



$$\tilde{m}_{\theta}(s_1, s_2) = F_{\theta_F}(s_1)^{\mathsf{T}} B_{\theta_B}(s_2)$$

**Definition 1** (Forward-backward representation). Let  $Z = \mathbb{R}^d$  be a representation space, and let  $\rho$ be a measure on  $S \times A$ . A pair of functions  $F: S \times A \times Z \to Z$  and  $B: S \times A \to Z$ , together with a parametric family of policies  $(\pi_z)_{z \in Z}$ , is called a forward-backward representation of the MDP with respect to  $\rho$ , if the following conditions hold for any  $z \in Z$  and  $(s, a), (s_0, a_0) \in S \times A$ :

$$\pi_z(s) = \arg\max_a F(s, a, z)^{\mathsf{T}} z, \qquad M^{\pi_z}(s_0, a_0, \mathrm{d}s, \mathrm{d}a) = F(s_0, a_0, z)^{\mathsf{T}} B(s, a) \rho(\mathrm{d}s, \mathrm{d}a)$$
(2)

where  $M^{\pi}$  is the successor measure defined in (1), and the last equality is between measures. **Theorem 2** (FB representations encode all optimal policies). Let  $(F, B, (\pi_z))$  be a forward-backward representation of a reward-free MDP with respect to some measure  $\rho$ .

Then, for any bounded reward function  $r: S \times A \to \mathbb{R}$ , the following holds. Set

$$z_R := \int_{s,a} r(s,a) B(s,a) \,\rho(\mathrm{d}s,\mathrm{d}a). \tag{3}$$

assuming the integral exists. Then  $\pi_{z_R}$  is an optimal policy for reward r in the MDP. Moreover, the optimal Q-function  $Q^*$  for reward r is  $Q^*(s, a) = F(s, a, z_R)^T z_R$ .

### **Matrix Factorization**

• It provides a direct representation of the value function at every state, without learning an additional model of V. Namely,

$$V(s) \approx F(s)^{\top} B(R), \qquad B(R) := \mathbb{E}_{s \sim \rho}[r_s B(s)]$$

$$(47)$$

where the "reward representation" B(R) can be directly estimated by an online average of B(s) weighted by the reward  $r_s$  at s. This is a direct consequence of (17). For instance, with sparse rewards, each time a reward is observed, the value function is updated everywhere. 10

This point applies to successor states, but not to goal-dependent value functions, which cannot handle arbitrary rewards.

• It simplifies the sampling of a pair of states  $(s, s_2)$  needed for forward TD. Indeed, the forward TD update (21) factorizes as an expectation over s, times an expectation over  $s_2$  (Section 6.2 below), which can be estimated independently. The same applies to backward TD. This can potentially reduce variance a lot, and even allows for purely "trajectorywise" online estimates using only the current transition  $s \to s'$ , without sampling of another independent state  $s_2$ . (Once more, this works for successor states but not for goal-dependent value functions, since in that case the transitions  $s \to s'$  depend on  $s_2$ .)



- It produces two (policy-dependent) representations of states, a forward and a backward one, in a natural way from the dynamics of the MDP and the current policy. This could be useful for other purposes.
- Even in the tabular case, when the state space is discrete and unstructured, this provides a form of prior or generalization between states (based on a low-rank prior for the successor state operator). States that are linked by the MDP dynamics get representations F and Bthat are close.
- It has some of the properties of the second-order methods of Section 7, without their complexity. This is proved in Appendix F.1.

The shortcomings are as follows:

- It approximates the successor state operator by an operator of rank at most r. This is never an exact representation unless the representation dimension r is at least the number of distinct states.
- The best rank-r approximation of (Id -γP)<sup>-1</sup> erases the small singular values of P: thus this representation will tend to erase "high frequencies" in the reward and value function, and provide a spatially smoother approximation focusing on long-range behavior. This is fine as long as the reward is not a "fast-changing" function made up of high frequencies (such as a "checkerboard" reward).

This can be expected: learning a reward-agnostic object such as M cannot work equally well for all rewards. For these reasons, it may be useful to use a mixed model for the value function with the FB model as a baseline, such as

$$V_{\varphi}(s) = F(s)^{\top} B(R) + v_{\varphi}(s) \tag{48}$$

where F and B are learned via successor states, B(R) is as in (47), and  $\varphi$  is learned via ordinary TD on the remainder. The FB part will catch reward-independent, long-range behavior, while the  $v_{\varphi}$  part will be needed to catch high frequencies in a particular reward function.

# Code



Algorithm 1 FB algorithm: Unsupervised Phase

```
1: Inputs: replay buffer \mathcal{D}, Polyak coefficient \alpha, \nu a probability distribution over \mathbb{R}^d, randomly
          initialized networks F_{\theta} and B_{\omega}, learning rate \eta, mini-batch size b, number of episodes E, number
          of gradient updates N, temperature \tau and regularization coefficient \lambda.
   2: for m = 1, ..., do
   3:
                /* Collect E episodes
   4:
                for episode e = 1, \ldots E do
                      Sample z \sim \nu
   5:
                      Observe an initial state s_0
   6:
                      for t = 1, ..., do
   7:
                            Select an action a_t according to some behaviour policy (e.g the \varepsilon-greedy with respect to
   8:
                             F_{\theta}(s_t, a, z)^{\top} z)
   9:
                             Observe next state s_{t+1}
                            Store transition (s_t, a_t, s_{t+1}) in the replay buffer \mathcal{D}
10:
11:
                      end for
                end for
12:
                /* Perform N stochastic gradient descent updates
13:
                for n = 1 \dots N do
14:
                      Sample a mini-batch of transitions \{(s_i, a_i, s_{i+1})\}_{i \in I} \subset \mathcal{D} of size |I| = b.
15:
                      Sample a mini-batch of target state-action pairs \{(s'_i, a'_i)\}_{i \in I} \subset \mathcal{D} of size |I| = b.
16:
                      Sample a mini-batch of \{z_i\}_{i \in I} \sim \nu of size |I| = b.
17:
                      Set \pi_{z_i}(\cdot \mid s_{i+1}) = \operatorname{softmax}(F_{\theta^-}(s_{i+1}, \cdot, z_i)^\top z_i / \tau)
18:
19:
                      \mathscr{L}(\theta,\omega) =
                       \frac{1}{2b^2} \sum_{i,j \in I^2} \left( F_{\theta}(s_i, a_i, z_i)^{\top} B_{\omega}(s'_j, a'_j) - \gamma \sum_{a \in A} \pi_{z_i}(a \mid s_{i+1}) \cdot F_{\theta^-}(s_{i+1}, a, z_i)^{\top} B_{\omega^-}(s'_j, a'_j) \right)^2 - \frac{1}{2b^2} \sum_{i,j \in I^2} \left( F_{\theta}(s_i, a_i, z_i)^{\top} B_{\omega}(s'_j, a'_j) - \gamma \sum_{a \in A} \pi_{z_i}(a \mid s_{i+1}) \cdot F_{\theta^-}(s_{i+1}, a, z_i)^{\top} B_{\omega^-}(s'_j, a'_j) \right)^2 - \frac{1}{2b^2} \sum_{i,j \in I^2} \left( F_{\theta}(s_i, a_i, z_i)^{\top} B_{\omega}(s'_j, a'_j) - \gamma \sum_{a \in A} \pi_{z_i}(a \mid s_{i+1}) \cdot F_{\theta^-}(s_{i+1}, a, z_i)^{\top} B_{\omega^-}(s'_j, a'_j) \right)^2 - \frac{1}{2b^2} \sum_{i,j \in I^2} \left( F_{\theta}(s_i, a_i, z_i)^{\top} B_{\omega}(s'_j, a'_j) - \gamma \sum_{a \in A} \pi_{z_i}(a \mid s_{i+1}) \cdot F_{\theta^-}(s_{i+1}, a, z_i)^{\top} B_{\omega^-}(s'_j, a'_j) \right)^2 - \frac{1}{2b^2} \sum_{i,j \in I^2} \left( F_{\theta}(s_i, a_i, z_i)^{\top} B_{\omega^-}(s'_j, a'_j) - \gamma \sum_{a \in A} \pi_{z_i}(a \mid s_{i+1}) \cdot F_{\theta^-}(s_{i+1}, a, z_i)^{\top} B_{\omega^-}(s'_j, a'_j) \right)^2 - \frac{1}{2b^2} \sum_{i,j \in I^2} \left( F_{\theta}(s_i, a_i, z_i)^{\top} B_{\omega^-}(s'_j, a'_j) - \gamma \sum_{a \in A} \pi_{z_i}(a \mid s_{i+1}) \cdot F_{\theta^-}(s_{i+1}, a, z_i)^{\top} B_{\omega^-}(s'_j, a'_j) \right)^2 - \frac{1}{2b^2} \sum_{i,j \in I^2} \left( F_{\theta}(s_i, a_i, z_i)^{\top} B_{\omega^-}(s'_j, a'_j) - \gamma \sum_{a \in A} \pi_{z_i}(a \mid s_{i+1}) \cdot F_{\theta^-}(s'_i, a'_j) \right)^2 \right)^2 
                     \begin{array}{l} \frac{1}{b} \sum_{i \in I} F_{\theta}(s_i, a_i, z_i)^{\top} B_{\omega}(s_i, a_i) \\ /^{*} \text{ Compute orthonormality regularization loss} \\ \mathscr{L}_{\text{reg}}(\omega) &= \frac{1}{b^2} \sum_{i,j \in I^2} B_{\omega}(s_i, a_i)^{\top} \text{stop-gradient}(B_{\omega}(s'_j, a'_j)) \end{array}
20:
21:
                      \mathtt{stop-gradient}(B_{\omega}(s_i, a_i)^{\top} B_{\omega}(s'_i, a'_i)) - \frac{1}{h} \sum_{i \in I} B_{\omega}(s_i, a_i)^{\top} \mathtt{stop-gradient}(B_{\omega}(s_i, a_i))
                      Update \theta \leftarrow \theta - \eta \nabla_{\theta} \mathscr{L}(\theta, \omega) and \omega \leftarrow \omega - \eta \nabla_{\omega} (\mathscr{L}(\theta, \omega) + \lambda \cdot \mathscr{L}_{reg}(\omega))
22:
23:
                end for
               /* Update target network parameters
24:
                \theta^- \leftarrow \alpha \theta^- + (1-\alpha)\theta
25:
               \omega^- \leftarrow \alpha \omega^- + (1 - \alpha) \omega
26:
27: end for
```

# Discussion



$$\psi^{\pi}(s_{0}, a_{0}) := \mathbb{E}\left[\sum_{t \ge 0} \gamma^{t} \varphi(s_{t+1}) \mid s_{0}, a_{0}, \pi\right].$$

$$Q^{\pi}(s, a) = \mathbb{E}^{\pi}\left[\sum_{i=t}^{\infty} \gamma^{i-t} \phi_{i+1} \mid S_{t} = s, A_{t} = a\right]^{\top} \mathbf{w} \equiv \psi^{\pi}(s, a)^{\top} \mathbf{w}, \qquad w = reward/\phi$$

$$\pi_{z}(s) = \operatorname*{arg\,max}_{a} F(s, a, z)^{\top} z, \qquad M^{\pi_{z}}(s_{0}, a_{0}, \mathrm{d}s, \mathrm{d}a) = F(s_{0}, a_{0}, z)^{\top} B(s, a) \rho(\mathrm{d}s, \mathrm{d}a) \qquad z_{R} := \mathbb{E}[r(s, a)B(s, a)]$$

If given the same representation for  $\Phi$  and B in finite space, the two Q values will be totally same! Where F is just a discount matrix to calculate successor features. It means FB =  $\varphi$ .

But the two values differs when it is just a unsupervised feature extraction. And the procedure of inference is totally different.

# Some experiments-SFs





Figure 4: 49 randomly sampled reward functions learned by VISR.



Figure 5: The approximations to the optimal value functions for the reward functions in Figure 4, computed by VISR through GPI on 10 randomly sampled policies.

# Some experiments-FB





Figure 3: Heatmap of  $\max_a F(s, a, z_R)^\top z_R$  for  $z_R = B(\bigstar)$  Left: d = 25. Right: d = 75.



Figure 5: Trajectories generated by the  $F^{\top}B$  policies for the task of reaching a target position (star shape rightarrow while avoiding forbidden positions (red shape  $\bullet$ )



Figure 4: Contour plot of  $\max_{a \in A} F(s, a, z_R)^{\top} z_R$  in Continuous Maze. Left: for the task of reaching a target while avoiding a forbidden region, **Right**: for two equally rewarding targets.



Figure 6: Trajectories generated by the  $F^{\top}B$  policies for the task of reaching the closest among two equally rewarding positions (star shapes  $\bigstar$ ). (Optimal *Q*-values are not linear over such mixtures.)

# Some experiments-FB



Figure 14: Continuous Maze: Contour plots plot of  $\max_{a \in A} F(s, a, z_R)^{\top} z_R$  (left) and trajectories of the  $\varepsilon$  greedy policy with respect to  $F(s, a, z_R)^{\top} z_R$  with  $\varepsilon = 0.1$  (right). Left: for the task of reaching a target while avoiding a forbidden region, Middle: for the task of reaching the closest goal among two equally rewarding positions, **Right**: choosing between a small, close reward and a large, distant one..





模式识别与神经计算研究组

PAttern Recognition and NEural Computing

ParN<sub>e</sub>C

Figure 15: Ms. Pacman: Trajectories of the  $\varepsilon$  greedy policy with respect to  $F(s, a, z_R)^{\top} z_R$  with  $\varepsilon = 0.1$  (right). Top row: for the task of reaching a target while avoiding a forbidden region, Middle row: for the task of reaching the closest goal among two equally rewarding positions, Bottom row: choosing between a small, close reward and a large, distant one..



Figure 18: Continuous maze: Visualization of FB embedding vectors after projecting them in two-dimensional space with t-SNE. Left: the states to be mapped. Middle: the F embedding. Right: the B embedding.

# Thanks