

## Setting and contributions

### Episodic RL

- Bilinear exponential family (BEF) model:

$$\begin{aligned}\mathbb{P}(\tilde{s} | s, a) &= \exp(\psi(\tilde{s})^\top M_{\theta^p} \varphi(s, a) - Z_{s,a}^p(\theta^p)) \\ \mathbb{P}(r | s, a) &= \exp(r B^\top M_{\theta^r} \varphi(s, a) - Z_{s,a}^r(\theta^r)).\end{aligned}$$

- Minimizes (pseudo-)regret:

$$\mathcal{R}(K) \triangleq \sum_{k=1}^K \left( V_{\theta,1}^{\pi^*}(s_1^k) - V_{\theta,1}^{\pi^t}(s_1^k) \right).$$

**Contributions:** We investigate episodic RL problem with unknown rewards and transition. Our contributions are:

- A **Linear value observation** for BEF transitions, **generalizing the Gaussian** transition observation of [1]
- An algorithm with: **tractable exploration**, **tractable planning**, and a  $\tilde{\mathcal{O}}(\sqrt{d^3 H^3 K})$  **regret upper-bound**.
- A **clipping-free** algorithm thanks to an **improved elliptical lemma**.

## BEF – RLSVI algorithm

BEF – RLSVI is similar to RLSVI, and is clipping-free.

### Algorithm 1 BEF – RLSVI

- Input:** failure rate  $\delta$ , constants  $\alpha^p, \eta$  and  $(x_k)_{k \in [K]} \in \mathbb{R}^+$
- for** episode  $k = 1, 2, \dots, \mathbf{do}$
- Observe initial state  $s_1^k$
- Sample noise  $\xi_k \sim \mathcal{N}(0, x_k (G^p)^{-1})$  such that

$$G^p = \frac{\eta}{\alpha^p} \mathbb{A} + \sum_{\tau=1}^{k-1} \sum_{h=1}^H (\varphi(s_h^\tau, a_h^\tau)^\top A_i^\top A_j \varphi(s_h^\tau, a_h^\tau))_{i,j \in [d]}$$

- Perturb reward parameter:  $\tilde{\theta}^r(k) = \hat{\theta}^r(k) + \xi_k$
- Compute  $(\tilde{Q}_h^k)_{h \in [H]}$  via Bellman-backtracking, see Algorithm 2
- for**  $h = 1, \dots, H$  **do**
- Pull action  $a_h^k = \arg \max_a \tilde{Q}_h^k(s_h^k, a)$
- Observe reward  $r(s_h^k, a_h^k)$  and state  $s_{h+1}^k$ .
- end for**
- Update the penalized ML estimators  $\hat{\theta}^p(k), \hat{\theta}^r(k)$
- end for**

Unlike optimistic approaches, exploration here is explicit and efficient as it does not involve a high-dimensional optimization.

### Algorithm 2 Bellman Backtracking

- Input** Parameters  $\hat{\theta}^p, \hat{\theta}^r$ , initialize  $\tilde{\theta} = (\tilde{\theta}^r, \hat{\theta}^p)$  and  $\forall s, V_{H+1}(s) = 0$
- for** steps  $h = H - 1, H - 2, \dots, 0$  **do**
- Calculate  $Q_{\tilde{\theta},h}(s, a) = \mathbb{E}_{s,a}^{\tilde{\theta}^r}[r] + \langle \phi^p(s, a), \int V_{\tilde{\theta},h+1}(s') \mu^p(s') ds' \rangle_{\mathcal{H}}$ .
- end for**

## Why is BEF – RLSVI tractable?

### Planning:

For an MDP of the BEF, we can write the state-action value function linearly, at step  $h$ :

$$\tilde{Q}_h^\pi(s, a) = \mathbb{E}^{\tilde{\theta}^r}[r(s, a)] + \left\langle \phi^p(s, a), \int_S \mu^p(\tilde{s}) \tilde{V}h + 1^\pi(\tilde{s}) d\tilde{s} \right\rangle.$$

### Key facts:

- $\phi^p$  and  $\psi^p$  are in an RKHS, i.e. **infinite dimensional**
- Using **Random Fourier Transform** entails  $\mathcal{O}(pH^2K \log(HK))$  dimensional approximations of  $\phi^p$  and  $\psi^p$
- Therefore, the planning has a  $\mathcal{O}(pH^3K \log(HK))$  **complexity**, pseudo-polynomial in  $p, H$  and  $K$ , thus tractable.

### Maximum likelihood estimation:

There are different methods to approximate ML estimator:

- Integral approximation** techniques:
  - Simulated annealing and importance sampling
  - MCMC techniques for approximating the partition function.
  - Optimizing a different objective, *the contrastive divergence*, yields a good approximation.
- If the **natural parameter and support** of the distribution are **bounded**, an  $\epsilon$ -approximation can be derived in  $\mathcal{O}(\text{poly}(k/\epsilon))$
- Score matching**: avoids approximating the partition function. Under certain conditions, the estimation can be solved in  $\mathcal{O}(d^3)$

## Regret bound

**Theorem (regret bound):** Let  $G_{s,a} \triangleq (\varphi(s, a)^\top A_i^\top A_j \varphi(s, a))_{i,j \in [d]}$  and  $\mathbb{A} \triangleq (\text{tr}(A_i A_j^\top))_{i,j \in [d]}$ . Under regularity of the Hessian and assuming

- $\max\{\|\theta^r\|_{\mathbb{A}}, \|\theta^p\|_{\mathbb{A}}\} \leq B_{\mathbb{A}}, \|\mathbb{A}^{-1} G_{s,a}\| \leq B_{\varphi, \mathbb{A}}$  and  $\mathbb{E}^{\theta^r}[r(s, a)] \in [0, 1]$
- The noise  $\xi_k \sim \mathcal{N}(0, x_k (G^p)^{-1})$  satisfies  $x_k \gtrsim dH^2$

then for all  $\delta \in (0, 1]$ , with probability at least  $1 - 7\delta$ ,

$$\mathcal{R}(K) = \mathcal{O}(\sqrt{d^3 H^3 K}).$$

### Tightness of regret upper-bound:

- A lower bound for episodic RL with continuous state-action spaces is still missing.
- For tabular RL, [2] proves a lower bound of order  $\Omega(\sqrt{H^3 SAK})$
- A tabular MDP is also a BEF model with  $d = S^2 \times A$
- BEF – RLSVI's yields  $R(K) = \mathcal{O}(\sqrt{(S^2 A)^3 H^3 K})$ , **tight in  $H$  and  $K$** .

## Interesting proof bits

**Optimism:** Key reasons for choosing RLSVI-type algorithms:

- Perturbing the reward** estimation guarantees **optimism with a constant probability**
- A constant probability of optimism is **enough to control the value function approximation error**

**Transportation:** Using **transportation inequalities** instead of the **simulation lemma** (c.f. Lemma 1 in [1]) reduces a  $\sqrt{H}$  regret factor

### Elliptical lemma:

- Leveraging the boundedness of the true value function enables using an **improved elliptical lemma** ( $\sqrt{H}$  less than [3])
- The **norm of features** can only be large  $\mathcal{O}(d)$  times, thus, we can **omit clipping** and reduce the regret by  $\sqrt{d}$  compared to [4].

### Approximate planning:

- To guarantee a tractable planning, we approximate the transition with  $(1/\sqrt{H^2 K})$ -error. Using **mis-specification style analysis**, we show that the approximation doesn't hinder the regret bound.
- Using a **Linear-RL algorithm directly** on top of the approximation would **lead to a linear regret**.

## Conclusion

For episodic RL with BEF transitions and rewards:

- We propose **BEF – RLSVI** that achieves a  $\mathcal{O}(\sqrt{d^3 H^3 K})$  regret
- We show that **tractable planning and exploration** are possible
- We give the second **example of continuous linear MDPs** in literature, although both are infinite dimensional

For linear RL style analyses: The **occurrences of values outside the plausible range**, e.g.  $\hat{V} \notin [0, H]$ , are finite

Future work:

- The paper could be complemented by **experimental evaluations** on relevant tasks.
- The **tractability of planning** can be extended to any **shift invariant kernel**: this can lead to interesting generalizations.

**Acknowledgements:** The authors acknowledge the funding of the French National Research Agency, the French Ministry of Higher Education and Research, Inria, the MEL and the I-Site ULNE regarding project R-PILOTE-19-004-APPRENF

## References

- [1] Tongzheng Ren, Tianjun Zhang, Csaba Szepesvári, and Bo Dai. A free lunch from the noise: Provable and practical exploration for representation learning. In *Uncertainty in Artificial Intelligence*. PMLR, 2022.
- [2] Omar Darwiche Domingues, Pierre Ménard, Emilie Kaufmann, and Michal Valko. Episodic reinforcement learning in finite mdps: Minimax lower bounds revisited. In *Algorithmic Learning Theory*. PMLR, 2021.
- [3] Sayak Ray Chowdhury, Aditya Gopalan, and Odalric-Ambrym Maillard. Reinforcement learning in parametric mdps with exponential families. In *AISTATS*. PMLR, 2021.
- [4] Andrea Zanette, David Brandfonbrener, Emma Brunskill, Matteo Pirota, and Alessandro Lazaric. Frequentist regret bounds for randomized least-squares value iteration. In *AISTATS*. PMLR, 2020.