# Leveraging Factored Action Spaces for Efficient Offline RL in Healthcare

Shengpu Tang, Maggie Makar, Michael W. Sjoding, Finale Doshi-Velez, Jenna Wiens
NeurIPS 2022

# Action Spaces in Clinical Problems

Commonly exhibit combinatorial structures

## Acute Dyspnea
**(ongoing project at UM)**

$|A| = 2^5 = 32$

💊 {0,1} Antibiotics

💊 {0,1} Anticoagulants

💉 {0,1} Fluids

💊 {0,1} Diuretics

💉 {0,1} Steroids

## Mech Vent Weaning
**(Prasad et al., UAI 2017)**

$|A| = 2 \times 4 = 8$

MV setting $\quad a[0] \in \{0, 1\}$

Sedation level $\quad a[1] \in \{0, 1, 2, 3\}$

$$\mathcal{A} = \left\{ \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 2 \end{bmatrix}, \begin{bmatrix} 0 \\ 3 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \begin{bmatrix} 1 \\ 3 \end{bmatrix} \right\}$$

## AI Clinician / MIMIC-sepsis
**(Komorowski et al., Nature Medicine 2018)**

$|A| = 5 \times 5 = 25$



|  | | Dose of vasopressor | | | |
|---|---|---|---|---|---|
|  | 1 | 2 | 3 | 4 | 5 |
| 1 | 1 | 2 | 3 | 4 | 5 |
| 2 | 6 | 7 | 8 | 9 | 10 |
| 3 | 11 | 12 | 13 | 14 | 15 |
| 4 | 16 | 17 | 18 | 19 | 20 |
| 5 | 21 | 22 | 23 | 24 | 25 |

Dose of i.v. fluid

$$\mathcal{A} = \mathcal{A}_1 \times \cdots \times \mathcal{A}_D$$

Overall action space is a
**Cartesian product** of $D$ sub-action spaces

$$\boldsymbol{a} = [a_1, \ldots, a_D] \in \mathcal{A}$$

$$a_d \in \mathcal{A}_d$$

Each action is a **vector** of
$D$ sub-actions

# Combinatorial action space → Typical Q function

$$Q(s, a)$$

Inefficient?

Tang et al., "Leveraging Factored Action Spaces for Efficient Offline RL in Healthcare", NeurIPS 2022.

5

# Factored action space → Linear Q decomposition



$$Q(s, \square\blacksquare\square) = q_1(s, \square) + q_2(s, \blacksquare) + q_3(s, \square)$$

**linear Q-function decomposition**

$$Q^\pi(s, \boldsymbol{a}) = \sum_{d=1}^{D} q_d(s, a_d)$$



We develop an approach for offline RL with **factored action spaces** by learning **linearly decomposable** Q-functions.

- Provide new **theoretical insights** on its applicability

- Conduct **empirical evaluations** in the context of offline RL for healthcare

**linear Q-function decomposition**

$$Q^\pi(s, \boldsymbol{a}) = \sum_{d=1}^{D} q_d(s, a_d)$$



## "When does it work?"

Does linear decomposition always exist? Will using linear decomposition introduce bias?

## Sufficient Conditions for Zero Bias         *...yet are not necessary*

*D* "parallel" MDPs  $\rightarrow$  implicitly factorized MDP via state abstractions

Outside the regime of theoretical guarantees --

Implication of linear approximation on bias, variance, and policy optimality

**Reduced Variance**    **Bias-Variance Trade-off**

The number of free parameters of tabular MDP

$$|\mathcal{S}||\mathcal{A}| = |\mathcal{S}|(\textstyle\prod_{d=1}^{D} |\mathcal{A}_d|) \quad \rightarrow \quad |\mathcal{S}|\left((\textstyle\sum_{d=1}^{D} |\mathcal{A}_d|) - D + 1\right)$$

**Bias $\nRightarrow$ Suboptimality**

e.g., when two sub-actions "reinforce" their independent effects

Demonstrate, with examples, how **domain knowledge** may be used to inform its **applicability** in real-world problems (e.g., healthcare, education)

**Action Space:** $\mathcal{A} = \mathcal{A}_{\text{abx}} \times \mathcal{A}_{\text{vaso}} \times \mathcal{A}_{\text{mv}}$   $|\mathcal{A}| = 2^3 = 8$



$\rho = 0.125$

Policy Value

Sample Size

Behavior policy takes the optimal action with probability $\rho$

**Proposed approach…**

Outperforms **baseline** for small sample sizes
Closely matches **baseline** for large sample sizes

**Action Space:**  $\mathcal{A} = \mathcal{A}_{\mathrm{abx}} \times \mathcal{A}_{\mathrm{vaso}} \times \mathcal{A}_{\mathrm{mv}}$    $|\mathcal{A}| = 2^3 = 8$



$\rho = 0.01$    $\rho = 0$

Behavior policy takes the optimal action **less than random**

**Proposed approach** <span style="color:red">better at inferring **underexplored** actions</span>

Problem setup based on Komorowski et al., "AI Clinician", *Nature Medicine* 2018.

**State Space**

Derived from 48 physiological signals



**Action Space**

$$\mathcal{A} = \mathcal{A}_{\text{vaso}} \times \mathcal{A}_{\text{fluid}}$$
$$|\mathcal{A}| = 5 \times 5 = 25$$



| Policy | Baseline BCQ | Factored BCQ | Clinician |
|--------|--------------|--------------|-----------|
| Test WIS | $90.44 \pm 2.44$ | $91.62 \pm 2.12$ | $90.29 \pm 0.51$ |
| Test ESS | $178.32 \pm 11.42$ | $178.32 \pm 11.96$ | 2894 |

**Better performance** at same effective sample size

Tang et al., "Leveraging Factored Action Spaces for Efficient Offline RL in Healthcare", NeurIPS 2022.

12

See paper for details



**Clinician**

| IV fluid dose (mL/4h) | 0 | 0.001-0.08 | 0.08-0.2 | 0.2-0.45 | >0.45 |
|---|---|---|---|---|---|
| >2L | 357 | 39 | 90 | 156 | 211 |
| 1L-2L | 1401 | 103 | 167 | 268 | 277 |
| 500mL-1L | 2984 | 162 | 179 | 296 | 273 |
| 1-500 | 17147 | 791 | 654 | 882 | 606 |
| 0 | 8491 | 106 | 48 | 83 | 75 |

Vasopressor dose (µg/kg/min)

**Baseline BCQ**

| IV fluid dose (mL/4h) | 0 | 0.001-0.08 | 0.08-0.2 | 0.2-0.45 | >0.45 |
|---|---|---|---|---|---|
| >2L | 13 | 0 | 0 | 1 | 119 |
| 1L-2L | 0 | 0 | 0 | 0 | 6 |
| 500mL-1L | 4 | 0 | 0 | 0 | 14 |
| 1-500 | 22355 | 936 | 38 | 2508 | 13 |
| 0 | 9839 | 0 | 0 | 0 | 0 |

Vasopressor dose (µg/kg/min)

**Factored BCQ**

| IV fluid dose (mL/4h) | 0 | 0.001-0.08 | 0.08-0.2 | 0.2-0.45 | >0.45 |
|---|---|---|---|---|---|
| >2L | 153 | 0 | 3 | 177 | 65 |
| 1L-2L | 0 | 0 | 0 | 0 | 10 |
| 500mL-1L | 1244 | 0 | 1 | 183 | 77 |
| 1-500 | 22467 | 34 | 238 | 1801 | 154 |
| 0 | 9186 | 0 | 0 | 52 | 1 |

Vasopressor dose (µg/kg/min)

For less frequently observed / underexplored treatment combinations

Proposed approach captures their effects better

Tang et al., "Leveraging Factored Action Spaces for Efficient Offline RL in Healthcare", NeurIPS 2022.

13

# Takeaways

We develop an approach for offline RL with **factored action spaces** by learning **linearly decomposable** Q-functions.

- Leverage domain knowledge when available
- Identify scenarios when approximation bias does not lead to suboptimal performance
- Could apply more broadly to help scale RL methods in other applications involving combinatorial action spaces

S. Tang     M. Makar     M.W. Sjoding     F.Doshi-Velez     J. Wiens

This work is funded by NSF and NIH-NLM.

https://github.com/MLD3/OfflineRL_FactoredActions

$$s \rightarrow \sum \rightarrow Q(s,a)$$