

Composing Value Functions in Reinforcement Learning

Benjamin van Niekerk*, Steven James*,
Adam Earle, Benjamin Rosman

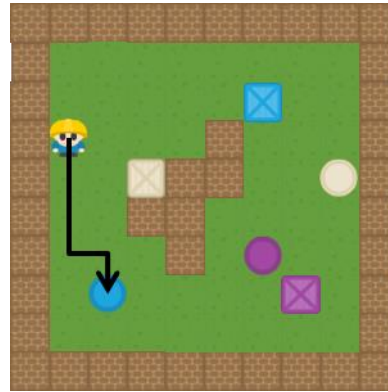
University of the Witwatersrand, Johannesburg, South Africa

Poster #251

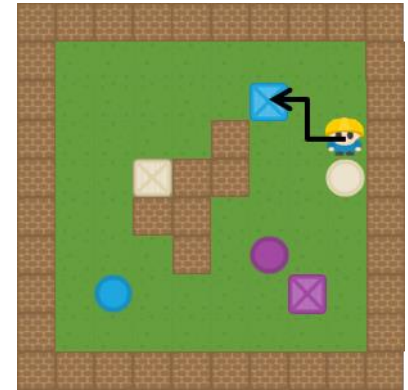
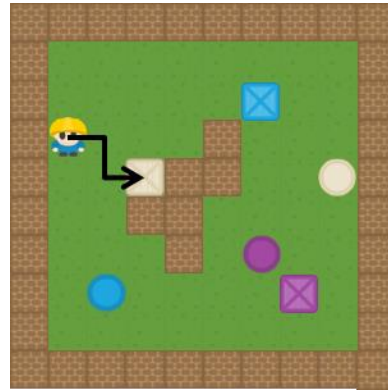
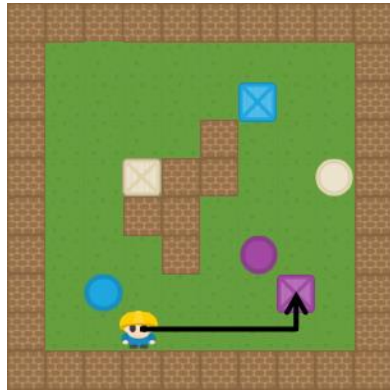
WITS
UNIVERSITY



$Q^*(\text{blue arrow})$:



$Q^*(\text{grey X})$:



Can we **blend** these value functions to solve interesting **combinations** of the tasks without further learning?

Skill Composition

- In general RL, for two skills Q_1 and Q_2 , we have that

Skill Composition

- In general RL, for two skills Q_1 and Q_2 , we have that

$$Q_1 \oplus Q_2 = \text{💩}$$

Entropy Regularised RL

- Augment reward function with penalty term:

$$r_{ent}(s, a) = r(s, a) - \tau \text{KL}[\pi_s \parallel \bar{\pi}_s]$$

$$V_{\pi}(s) = \mathbb{E}_s^{\pi} \left[\sum_{t=0}^{\infty} r(s_t, a_t) - \tau \text{KL}[\pi_{s_t} \parallel \bar{\pi}_{s_t}] \right]$$

Entropy Regularised RL

- Augment reward function with penalty term:

$$r_{ent}(s, a) = r(s, a) - \tau \text{KL}[\pi_s \parallel \bar{\pi}_s]$$

*Being different to
"reference" policy $\bar{\pi}$*

$$V_{\pi}(s) = \mathbb{E}_{\pi_s} \left[\sum_{t=0}^{\infty} r(s_t, a_t) - \tau \text{KL}[\pi_{s_t} \parallel \bar{\pi}_{s_t}] \right]$$

Entropy Regularised RL

- Augment reward function with penalty term:

$$r_{ent}(s, a) = r(s, a) - \tau \text{KL}[\pi_s \parallel \bar{\pi}_s]$$

Temperature

*Being different to
"reference" policy $\bar{\pi}$*

$$V_{\pi}(s) = \mathbb{E}_{\pi_s} \left[\sum_{t=0}^{\infty} r(s_t, a_t) - \tau \text{KL}[\pi_{s_t} \parallel \bar{\pi}_{s_t}] \right]$$

Skill Composition

- In general RL, for two skills Q_1 and Q_2 , we have that

$$Q_1 \oplus Q_2 = \text{💩}$$

Skill Composition

- In general RL, for two skills Q_1 and Q_2 , we have that

$$Q_1 \oplus Q_2 = \text{💩}$$

- With entropy regularisation, we show that

$$Q_1 \oplus Q_2 = \text{😍}$$

Provably ←

Skill Composition

- In general RL, for two skills Q_1 and Q_2 , we have that

$$Q_1 \oplus Q_2 = \text{💩}$$

- With entropy regularisation, we show that

$$Q_1 \oplus Q_2 = \text{😍}$$

Provably ←

*Deterministic
transitions only*

-OR- Task Composition

We can optimally compose $Q(\text{📦})$ and $Q(\text{🟪})$ to solve the task of collecting 📦 or 🟪 :

$$Q(\text{📦 or 🟪}) =$$

$$\tau \log(w_{\text{🟪}} \exp(Q(\text{🟪})/\tau) + w_{\text{📦}} \exp(Q(\text{📦})/\tau))$$

-OR- Task Composition

We can optimally compose $Q(\text{👛})$ and $Q(\text{🟪})$ to solve the task of collecting 👛 or 🟪 :

$$Q(\text{👛 or 🟪}) =$$

LogSumExp = "soft" maximum

$$\tau \log(w_{\text{🟪}} \exp(Q(\text{🟪})/\tau) + w_{\text{👛}} \exp(Q(\text{👛})/\tau))$$

-OR- Task Composition

We can optimally compose $Q(\text{👛})$ and $Q(\text{🟪})$ to solve the task of collecting 👛 or 🟪 :

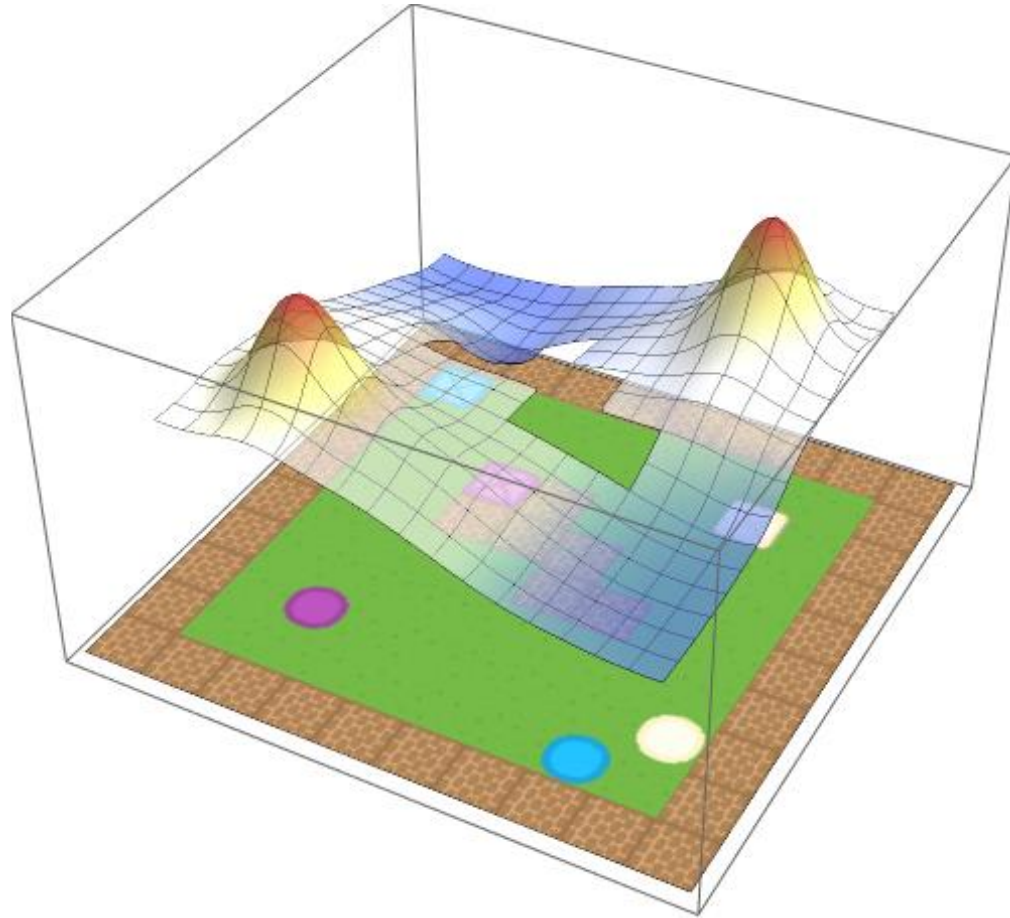
$$Q(\text{👛 or } \text{🟪}) =$$

LogSumExp = "soft" maximum

$$\tau \log(w_{\text{🟪}} \exp(Q(\text{🟪})/\tau) + w_{\text{👛}} \exp(Q(\text{👛})/\tau))$$

Use these to weight

Experiment -OR-



A Connection to Standard RL

Corollary: In the limit as $\tau \downarrow 0$ i.e. in the standard RL setting, we prove that:

$$Q(\text{ } \text{ or } \text{ }) = \max\{Q(\text{ }), Q(\text{ })\}$$

A Connection to Standard RL

Corollary: In the limit as $\tau \downarrow 0$ i.e. in the standard RL setting, we prove that:

$$Q(\text{ } \text{ or } \text{ }) = \max\{Q(\text{ }), Q(\text{ })\}$$

regular max

A Connection to Standard RL

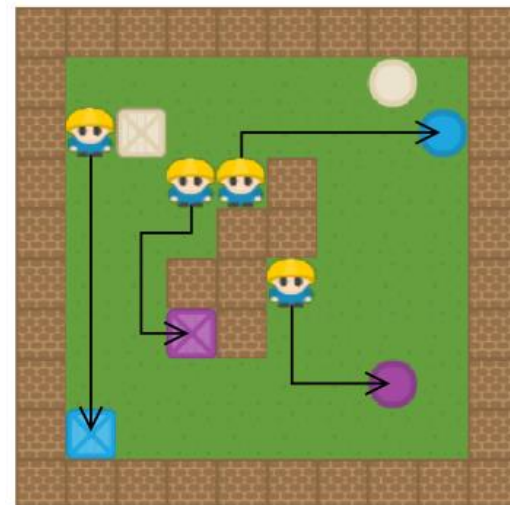
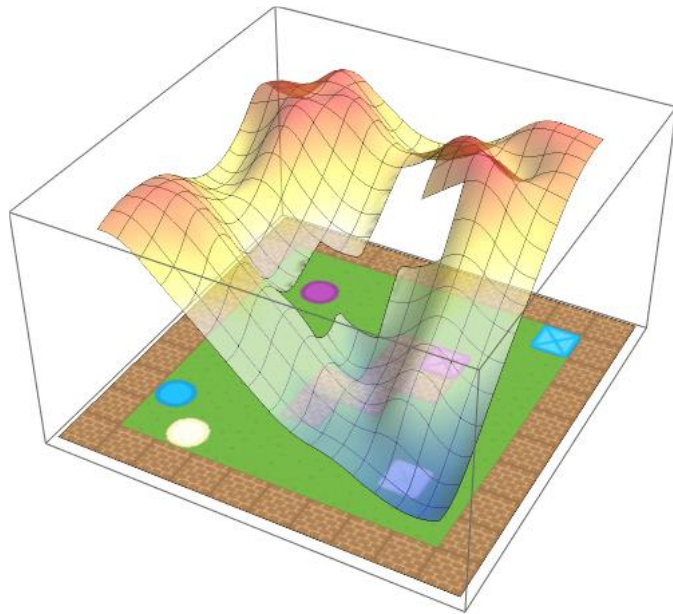
Corollary: In the limit as $\tau \downarrow 0$ i.e. in the standard RL setting, we prove that:

$$Q(\text{blue} \text{ or } \text{grey}) = \max\{Q(\text{blue}), Q(\text{grey})\}$$

regular max

No weights! ☹️

Experiment -OR-



-AND- Task Composition

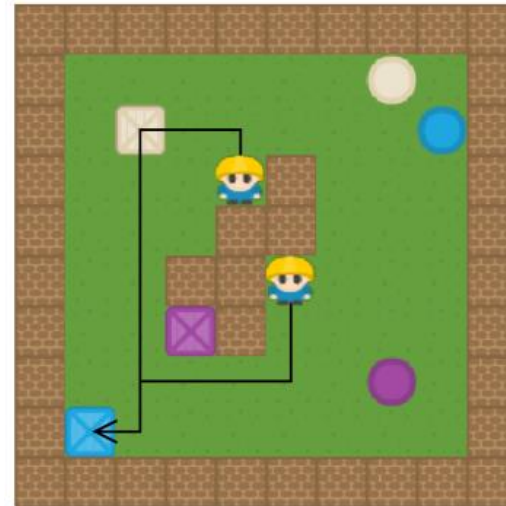
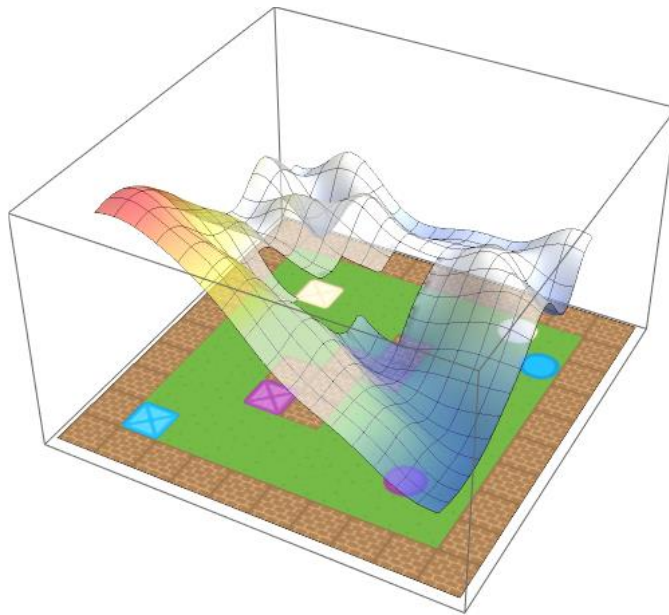
- Previous work in entropy regularised RL [Haarnoja et al, 2018] shows that -AND- task is approximately

$$Q(\text{👉 and } \text{☒}) \approx (Q(\text{👉}) + Q(\text{☒}))/2$$

- Conjecture this holds in low-temperature limit

Haarnoja, T., Pong, V., Zhou, A., Dalal, M., Abbeel, P., and Levine, S.
Composable deep reinforcement learning for robotic manipulation. *arXiv preprint arXiv:1803.06773, 2018.*

Experiment -AND-



Conclusion

- We can do zero-shot composition to provably find

$$Q^*(\text{👉 or } \text{☒}) \text{ (with different priorities)}$$

- We provide empirical evidence that averaging Q-values approximates:

$$Q^*(\text{👉 and } \text{☒})$$

Poster #251

