



Reinforcement Learning with Multiple Experts: A Bayesian Model Combination Approach

0712214 陳彥儒 0716007 潘冠蓁
0716308 張千祐 0716312 葉佳翰

Outline

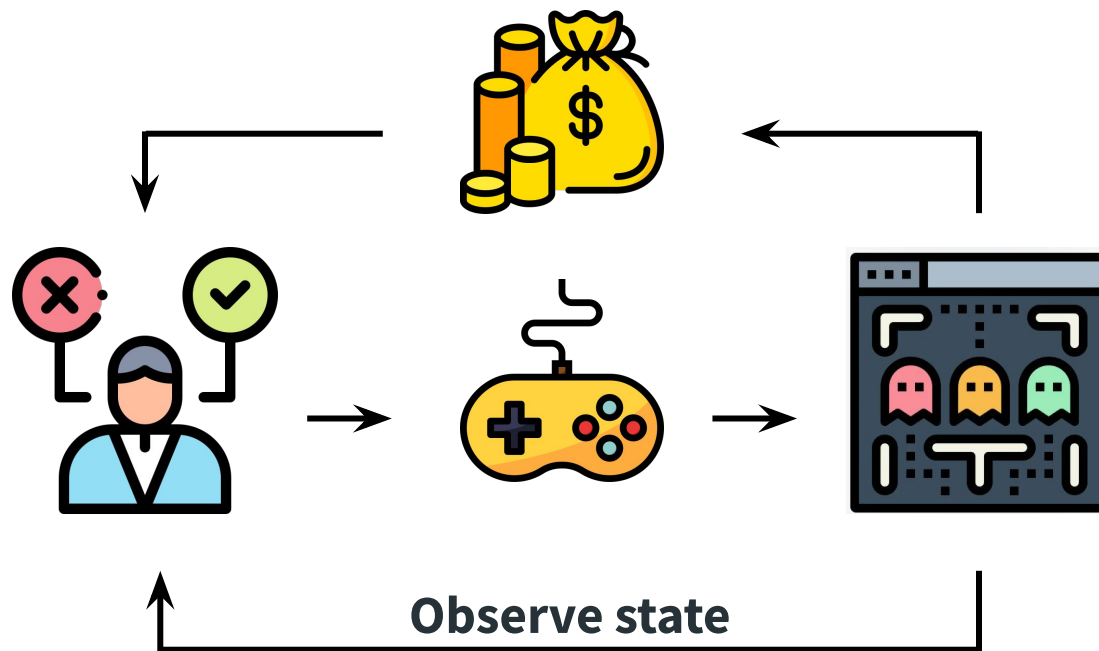
- ◎ Introduction
- ◎ Bayesian Reward Shaping
- ◎ Experimental Setting — Environment
- ◎ Experimental Setting — RL algorithm
- ◎ Experimental Setting — Experts
- ◎ Experimental Results
- ◎ Conclusion

A decorative network diagram in the top-left corner, featuring a complex web of interconnected nodes and lines. The nodes are represented by small circles, some of which are larger and have concentric circles, suggesting different levels of connectivity or importance. The lines are thin and gray, creating a mesh-like structure.

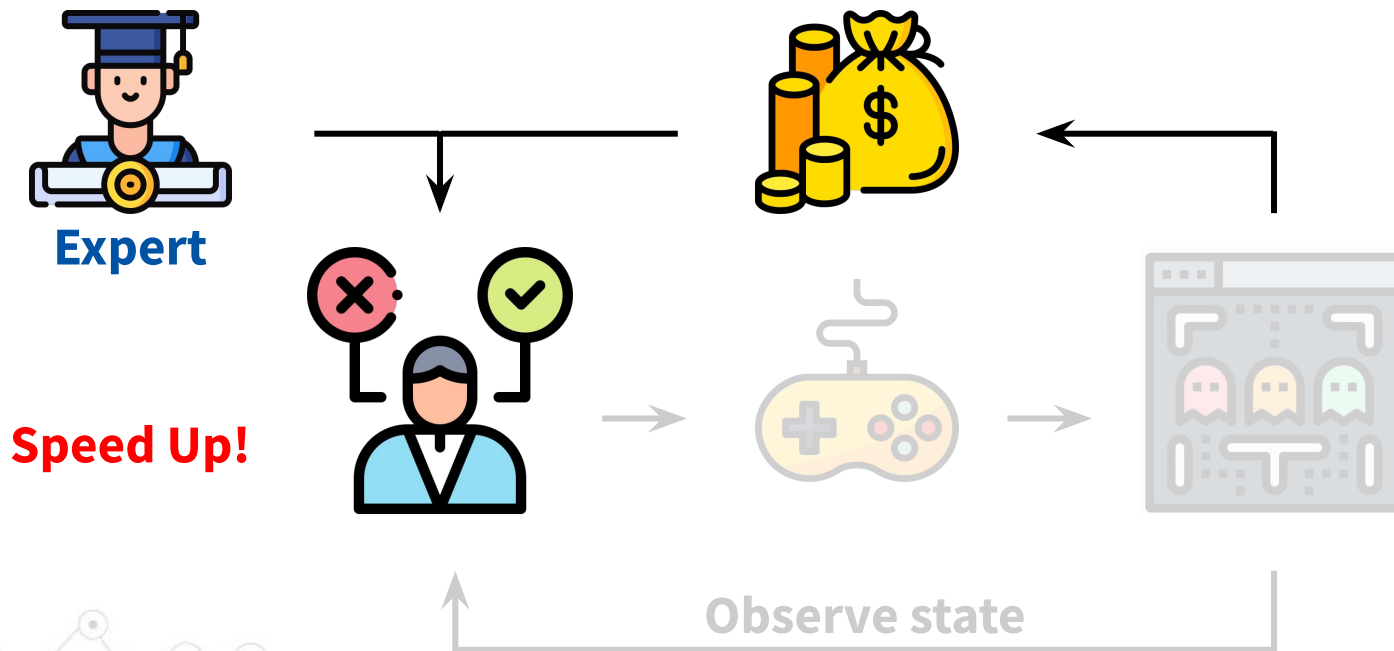
Introduction

A decorative network diagram in the bottom-right corner, similar to the one in the top-left. It shows a cluster of nodes connected by lines, with some nodes being larger and more prominent than others. The overall style is clean and modern, with a focus on geometric patterns.

Introduction

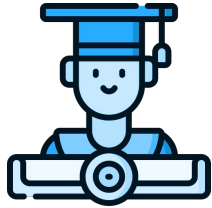


Introduction

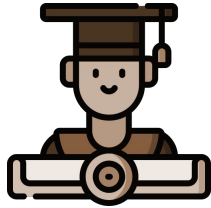


Introduction

But, if you have more than 1 experts...



Expert 1



Expert 2



Expert 3



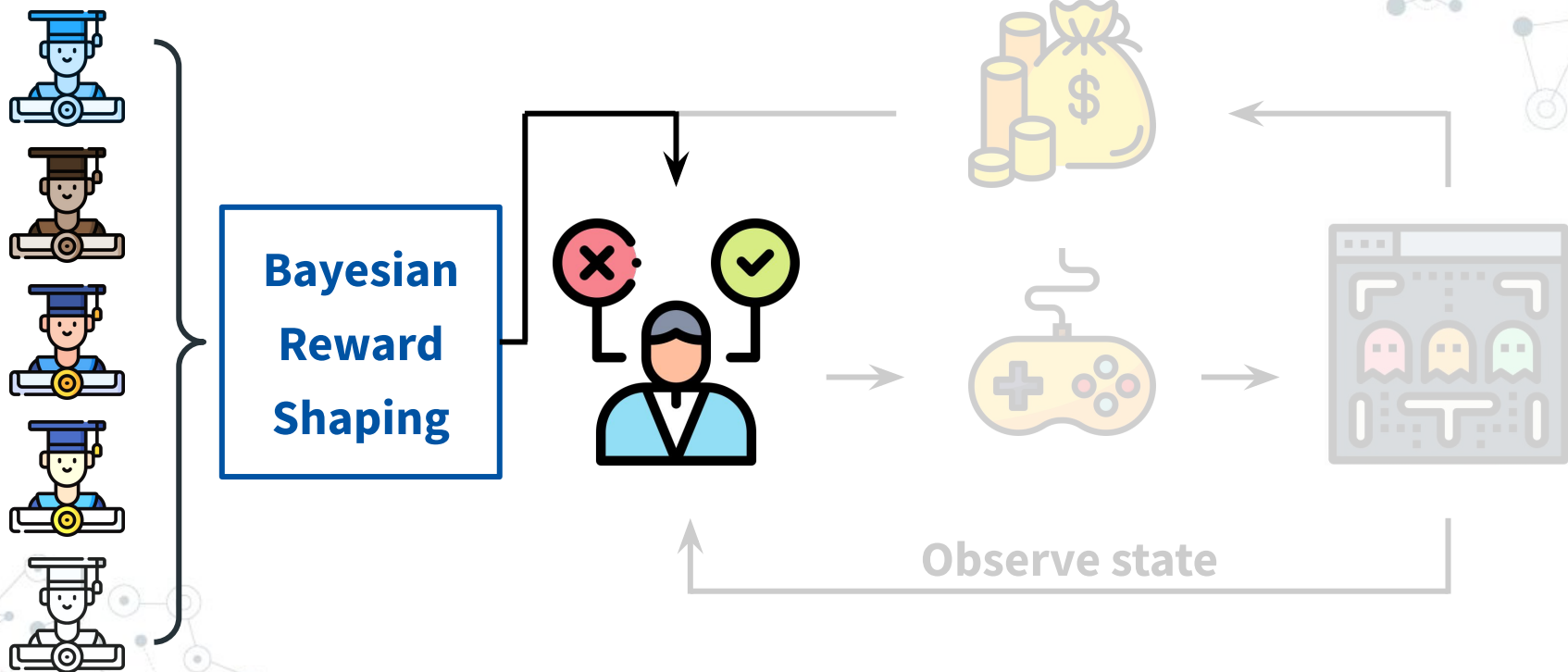
Expert 4



Expert 5

Which one should you trust?

Introduction





Bayesian Reward Shaping



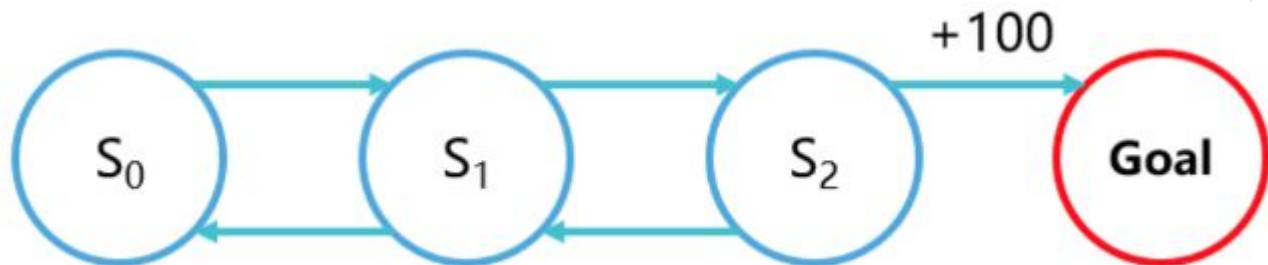
Diagram illustrating a state space search problem with states S_0 , S_1 , S_2 , and Goal. Transitions are shown between S_0 and S_1 , S_1 and S_2 , and S_2 and Goal. A reward of +100 is associated with the transition to the Goal state.

```

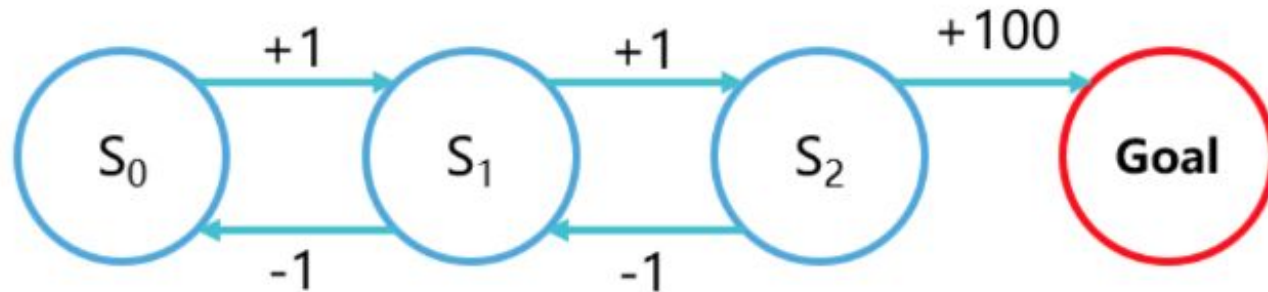
graph LR
    S0((S0)) -- +1 --> S1((S1))
    S1 -- +1 --> S2((S2))
    S2 -- +100 --> Goal((Goal))
    S0 -- +1 --> S0
  
```

Bayesian Reward Shaping

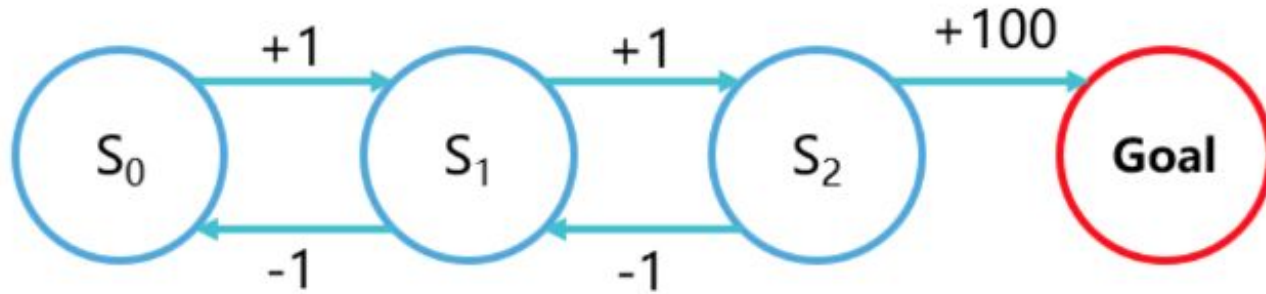
**Original
reward**



**But if we change a
little bit...**



Bayesian Reward Shaping



$$\Phi(s_0) = 0 \quad \Phi(s_1) = 1 \quad \Phi(s_2) = 2$$

$$\boxed{F(s, a, s')} = \gamma \Phi(s') - \Phi(s)$$

This might be the proper definition of
the additional reward from expert.

Bayesian Reward Shaping



$$\Phi_1 \longrightarrow F_1(s, s') = \gamma\Phi_1(s') - \Phi_1(s)$$



$$\Phi_2 \longrightarrow F_2(s, s') = \gamma\Phi_2(s') - \Phi_2(s)$$



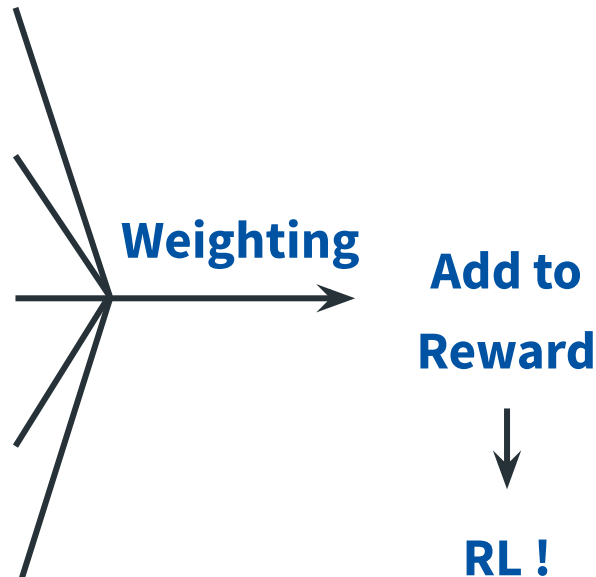
$$\Phi_3 \longrightarrow F_3(s, s') = \gamma\Phi_3(s') - \Phi_3(s)$$



$$\Phi_4 \longrightarrow F_4(s, s') = \gamma\Phi_4(s') - \Phi_4(s)$$



$$\Phi_5 \longrightarrow F_5(s, s') = \gamma\Phi_5(s') - \Phi_5(s)$$



Bayesian Reward Shaping

Algorithm 2 RL with Bayesian Reward Shaping

```
1: initialize  $\alpha \in \mathbb{R}_+^N$ 
2: for  $episode = 0, 1 \dots M$  do
3:    $\hat{\Phi} \leftarrow \frac{\sum_{i=1}^N \Phi_i \alpha_i}{\sum_{i=1}^N \alpha_i}$ 
4:    $F(s, a, s') \leftarrow \gamma \hat{\Phi}(s') - \hat{\Phi}(s)$ 
5:    $(R_t, s_t)_{t=1 \dots T} \leftarrow \text{TrainRL}(F)$ 
6:   for all  $(R_t, s_t)$  do
7:     update  $\hat{\sigma}^2$  and compute  $e$ 
8:      $\alpha \leftarrow \text{PosteriorUpdate}(\alpha, e)$ 
```

▷ Main loop

▷ Pool experts and compute shaped reward

▷ Perform one episode of training

▷ Posterior update

Bayesian Reward Shaping

Algorithm 2 RL with Bayesian Reward Shaping

```
1: initialize  $\alpha \in \mathbb{R}_+^N$ 
2: for  $episode = 0, 1 \dots M$  do                                ▷ Main loop
3:    $\hat{\Phi} \leftarrow \frac{\sum_{i=1}^N \Phi_i \alpha_i}{\sum_{i=1}^N \alpha_i}$           ▷ Pool experts and compute shaped reward
4:    $F(s, a, s') \leftarrow \gamma \hat{\Phi}(s') - \hat{\Phi}(s)$ 
5:    $(R_t, s_t)_{t=1 \dots T} \leftarrow \text{TrainRL}(F)$                 ▷ Perform one episode of training
6:   for all  $(R_t, s_t)$  do                                       ▷ Posterior update
7:     update  $\hat{\sigma}^2$  and compute  $e$ 
8:      $\alpha \leftarrow \text{PosteriorUpdate}(\alpha, e)$ 
```

How can we update the weights for each expert?

Bayesian Reward Shaping



Answer: lots of math...

Dirichlet Distribution: 估計 w 的分布
Target: Model how "properties" vary.
parameter: $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_k)$

$$P(w) \sim D(\alpha) = \frac{T(\sum_k \alpha_k)}{\prod_k T(\alpha_k)} \prod_k w_k^{\alpha_k - 1}, \quad \sum_k w_k = 1, w_k > 0$$

weight of experts

Hence $E[w_k] = \frac{\alpha_k}{\sum_k \alpha_k}$ Seems like Property of

$$E[w_k^2] = E[w_k] \frac{1 + \alpha_k}{1 + \sum_k \alpha_k}$$

We mentioned in paper Equation (14)

$$\Rightarrow E[w_k^2] = \frac{\alpha_k(\alpha_k + 1)}{\alpha_k(\alpha_k + 1) + 1}$$

The "second moment" S_k mentioned in paper Equation (15)

$$\begin{aligned} \pi(w) &= \frac{P(w, D, d)}{P(D, d)} \\ &= \frac{P(w, D, d)}{P(D)P(d)} \quad \text{D, d are indep (Mg given)} \\ &= \frac{P(w, D|d)}{P(D)} \quad \text{Def of conditional Prob} \\ &= \frac{P(w, D)}{P(D)} \quad \text{Bayes Theorem} \\ &= \frac{P(w, D)}{P(D)} \quad \text{Def of conditional Prob} \\ &= P(w|D)P(D) \quad \text{D, d are indep (Mg given)} \\ &= \pi_k(w)P(D) \\ &= \frac{1}{C_{k+1}} \sum_{d=1}^D P(d|z)P(e|w)\pi_k(w) \quad \text{normalization since } \sum \pi_{k+1} = 1 \\ &= \frac{1}{C_{k+1}} \sum_{d=1}^D P(d|z)P(e|w)\pi_k(w) \\ C_{k+1} &= \int_{S^{k+1}} \sum_{d=1}^D P(d|z)P(e|w)\pi_k(w) dw \\ &= \sum_{d=1}^D e_k E[w_k] \quad \text{P(d|z)} \end{aligned}$$

$$\begin{aligned} E[g_k | D] &= \int_R g P(g|D) dg \quad \text{all possible } g \text{-value} \\ &= \int_R g \int_{S^{k+1}} P(g, w|D) dw dg \quad \text{all possible combination of weight} \\ &= \int_R g \int_{S^{k+1}} \frac{P(g, w, D)}{P(D)} \frac{P(w, D)}{P(w, D)} dw dg \\ &= \int_R g \int_{S^{k+1}} \frac{P(g|w, D)}{\pi_k(w)} \frac{P(w|D)}{\pi_k(w)} dw dg \quad \text{The probability that weight } w \text{ emerge} \\ &= \int_R g \int_{S^{k+1}} \frac{P(g|z)}{\pi_k(w)} w_k \pi_k(w) dw dg \quad \text{why if } w \rightarrow \text{move to outside} \quad \text{Def of } E \\ &= \sum_{k=1}^K \int_R g P(g|z) \int_{S^{k+1}} w_k \pi_k(w) dw dg \quad \text{why if } w \rightarrow \text{move to outside} \quad \text{Def of } E \\ &= \sum_{k=1}^K \int_R g P(g|z) \frac{E[w_k]}{\pi_k} dg \quad \text{Def of } E \\ &= \sum_{k=1}^K E[w_k] \int_R g P(g|z) dg \quad \text{Def of } E \\ &= \sum_{k=1}^K E[w_k] E[g|z] \quad \text{Def of } E \end{aligned}$$

Conclusion: We can compute $g(s, a)$ by computing $E[w_k]$ & $E[g|z]$ "separately"

If you are interested:

<https://docs.google.com/document/d/1juUxpze40b18kQzbZft8ML3voZ-Fa3aXOyj0Cwo7VfK/edit?usp=sharing>

Bayesian Reward Shaping

Don't worry. Let's take it easier...

Algorithm 1 PosteriorUpdate(α_t (e))

1: **for** $i = 1, 2 \dots N - 1$ **do**

2: $m_i \leftarrow \frac{\alpha_{t,i}(e_i + e \cdot \alpha_t)}{(e \cdot \alpha_t)(\alpha_{t,0} + 1)}$

▷ Compute posterior moments

3: $s_1 \leftarrow \frac{\alpha_{t,1}(\alpha_{t,1} + 1)(2e_1 + e \cdot \alpha_t)}{(e \cdot \alpha_t)(\alpha_{t,0} + 1)(\alpha_{t,0} + 2)}$

4: $\alpha_{t+1,0} \leftarrow \frac{m_1 - s_1}{s_1 - m_1^2}$

▷ Compute α_{t+1}

5: **for** $i = 1, 2 \dots N - 1$ **do**

6: $\alpha_{t+1,i} \leftarrow m_i \alpha_{t+1,0}$

7: $\alpha_{t+1,N} \leftarrow \alpha_{t+1,0} - \sum_{i=1}^{N-1} \alpha_{t+1,i}$

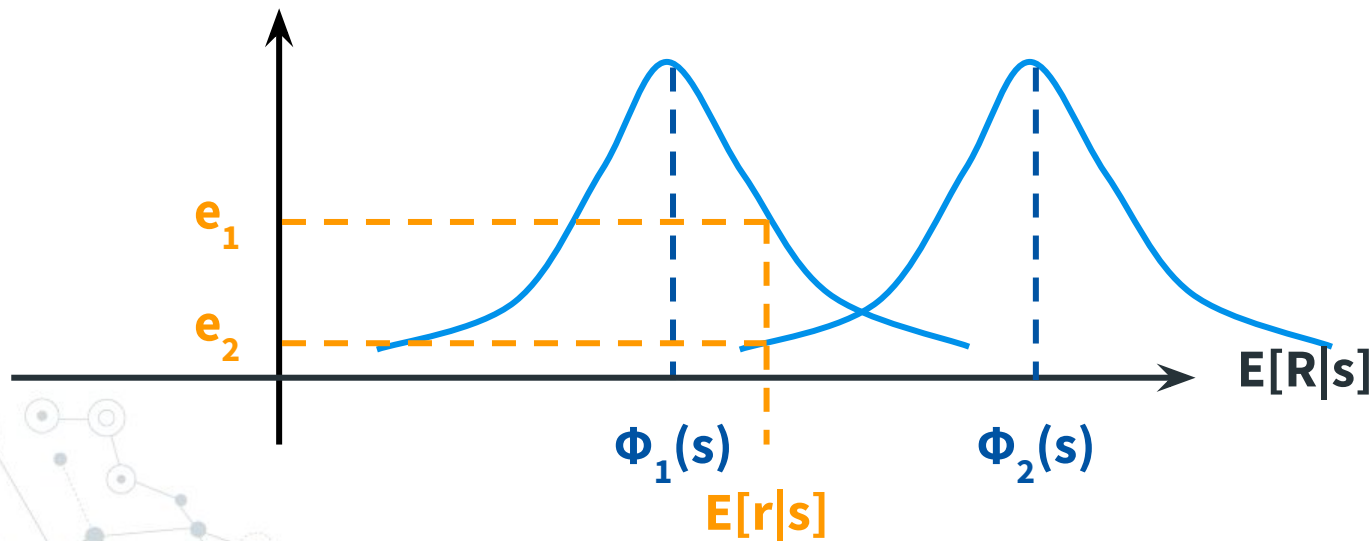
8: **return** α_{t+1}

“e” decide how to update the weights.

Bayesian Reward Shaping

What's “e” = “(e_1, e_2, \dots, e_n)”?

$e_i = P(E[r|s] \text{ occurs at expert } i)$

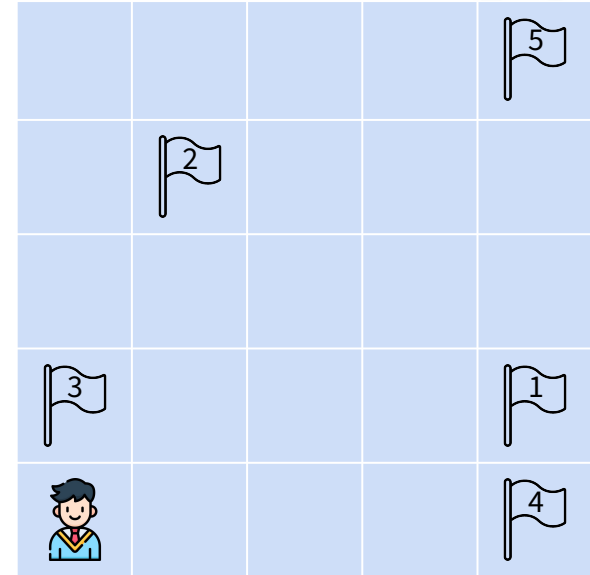


A decorative network diagram in the top-left corner, featuring a complex web of interconnected nodes and lines. The nodes are represented by small circles, some of which are larger and have concentric circles, suggesting different levels of connectivity or importance. The lines are thin and gray, creating a mesh-like structure.

Environment

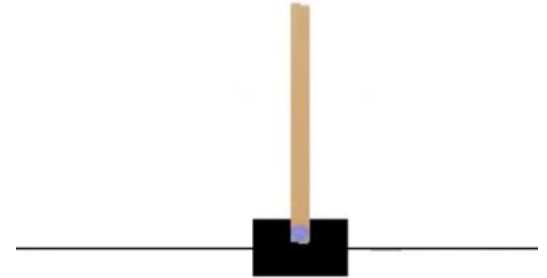
Gridworld

- ⦿ Rules
 - Every move : -1 point
 - Invalid move : -1 additional point
- ⦿ How to end this game
 - Until 200 steps
 - Collect all flags in order



Cartpole

- ◎ Goal
 - Keep the cartpole balanced
- ◎ When the game will end
 - Until 500 steps
 - $|\text{The angle between cart and pole}| > 12 \text{ deg}$
 - $|\text{The position of the cart}| > 2.4 \text{ deg}$

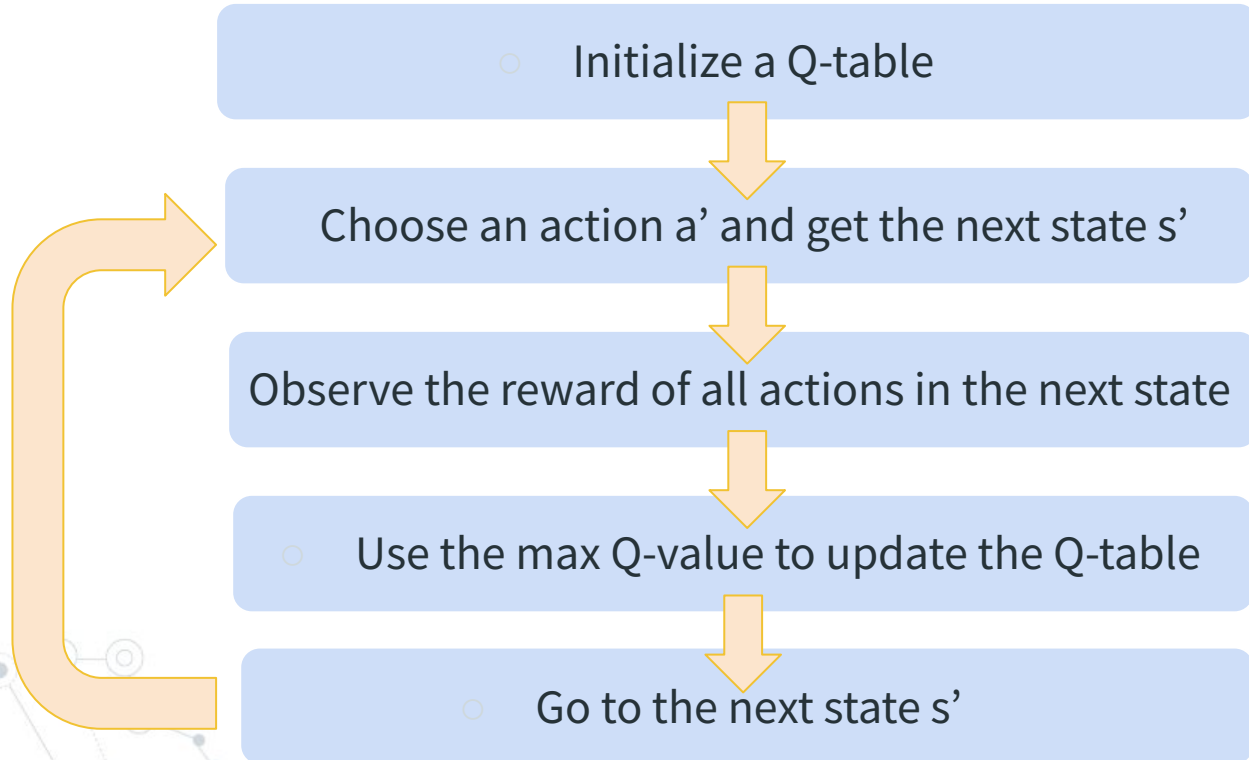




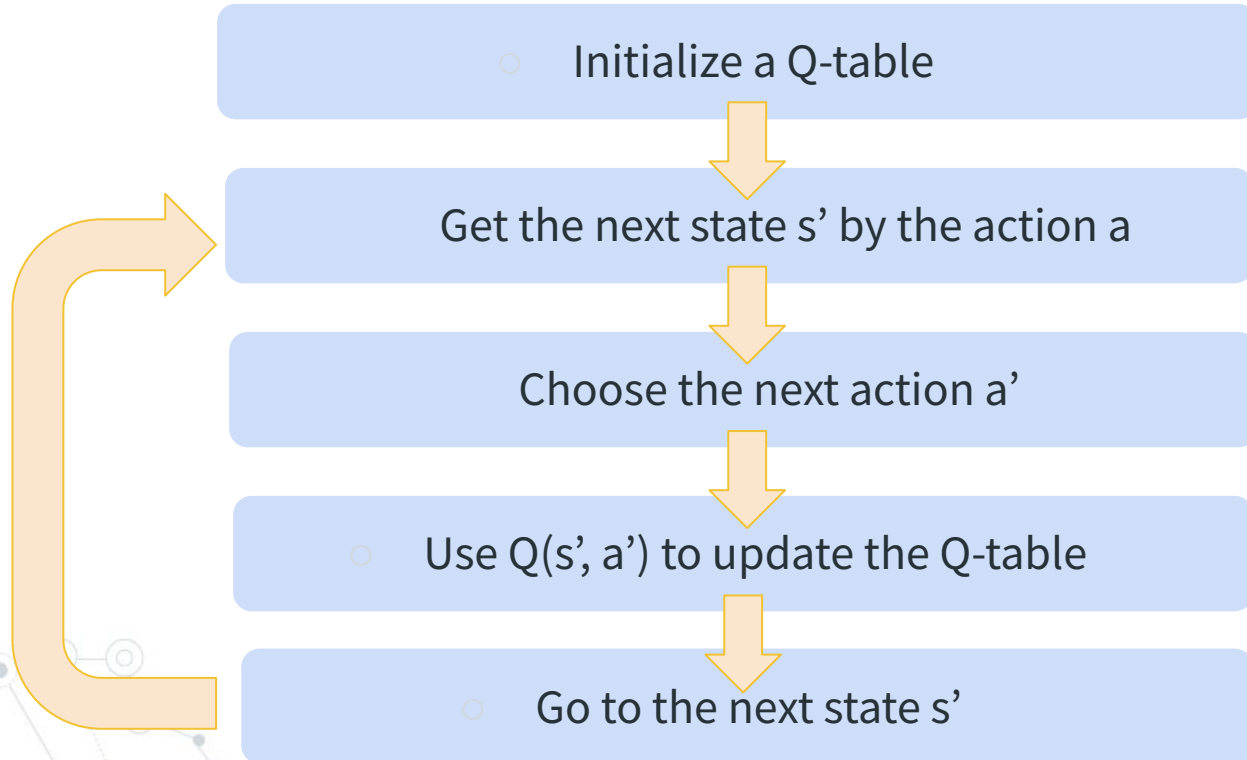
RL Algorithm



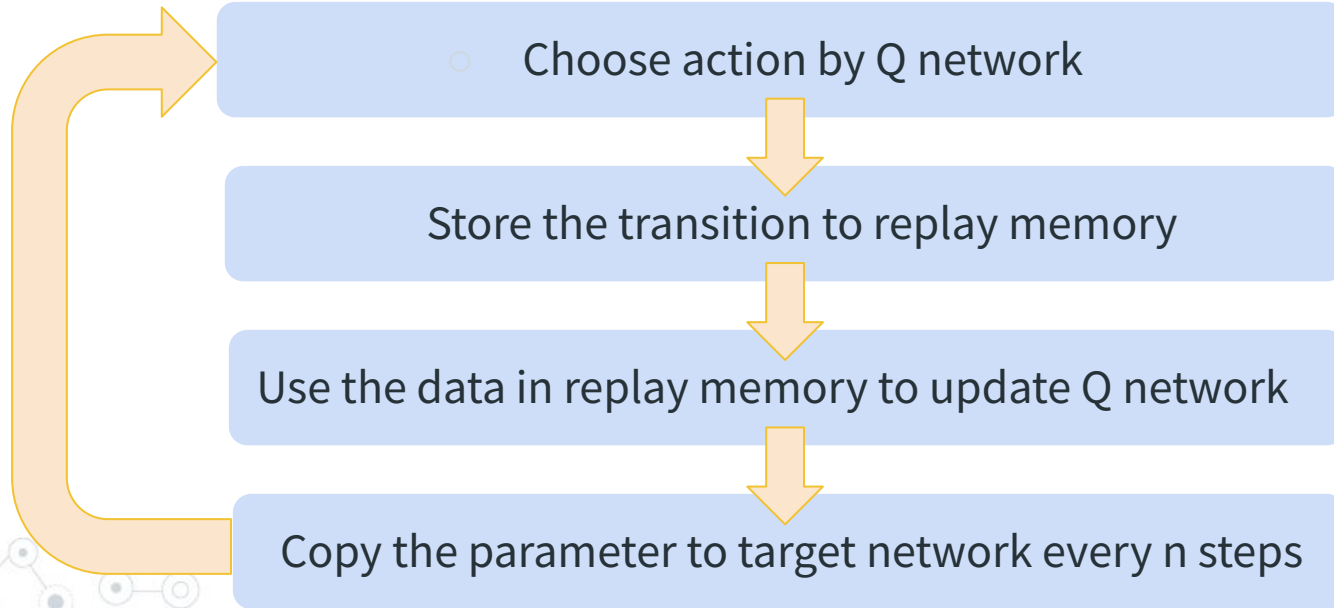
Tabular Q-learning



SARSA



Deep Q-Learning(DQN)



A decorative network diagram in the top-left corner, featuring a complex web of interconnected nodes and lines. The nodes are represented by small circles, some of which are larger and have concentric circles, suggesting different levels of connectivity or importance. The lines are thin and gray, creating a mesh-like structure.

Experts

Experts

For Gridworld:

1. $\Phi_{\text{good}}(s)$ = optimal value function
2. $\Phi_{\text{bad}}(s)$ = -optimal value function
3. $\Phi_{\text{random}}(s) = U(-20,20)$
4. $\Phi_{\text{zero}}(s) = 0$
5. $\Phi_{\text{heuristic}}(s) = -22 * (5 - c - 0.5) / 5$

Optimal value function:

$$V(x, y, c) = -\text{distance}((x, y), \text{next flag}) \\ - \text{distance}(\text{next flag}, \text{final})$$

Experts

For CartPole:

1. $\Phi_{\text{good}}(s) = Q$ Network trained by us

$$\Phi_{\text{good}}(s) = \text{avg } Q(s,a)$$

2. $\Phi_{\text{bad}}(s) = -Q$ Network trained by us

3. $\Phi_{\text{random}}(s) = U(-20,20)$

4. $\Phi_{\text{zero}}(s) = 0$

5. $\Phi_{\text{guess}}(s) = 20 * (1 - |\Theta| / 0.2618)$



Experimental Results



Experimental Results

Hyperparameters — Gridworld — DQN / Q-learning / SARSA

Learning rate: 0.001 / 0.4 / 0.36

Gamma: 0.99 / 1 / 1

Epsilon: $\max(0.98^t, 0.01)$ / 0.98^t / 0.98^t

Episode: 500 / 1000 / 2000

Independent set: 20 / 20 / 20

Experimental Results

Hyperparameters — CartPole — DQN / Q-learning / SARSA

Learning rate: 0.0005 / $\max(0.5 \cdot \text{lr}_{dc}^t, 0.01)$ / $\max(0.5 \cdot \text{lr}_{dc}^t, 0.01)$

Gamma: 0.99 / 0.95 / 0.95

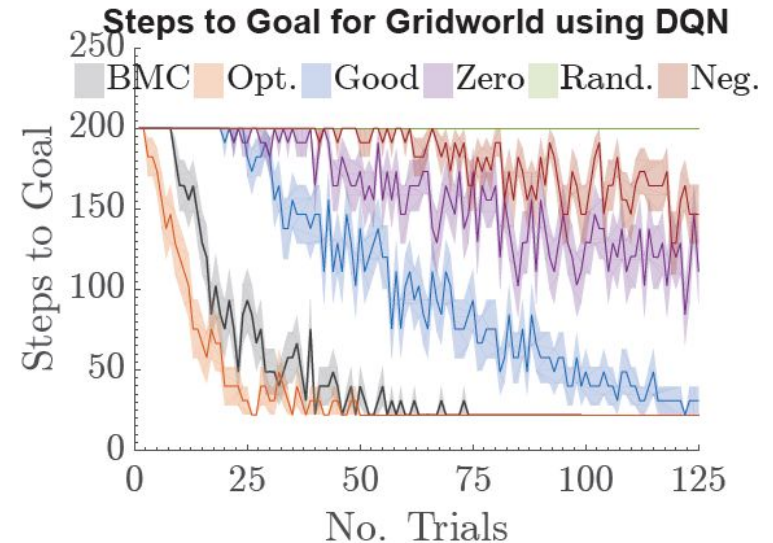
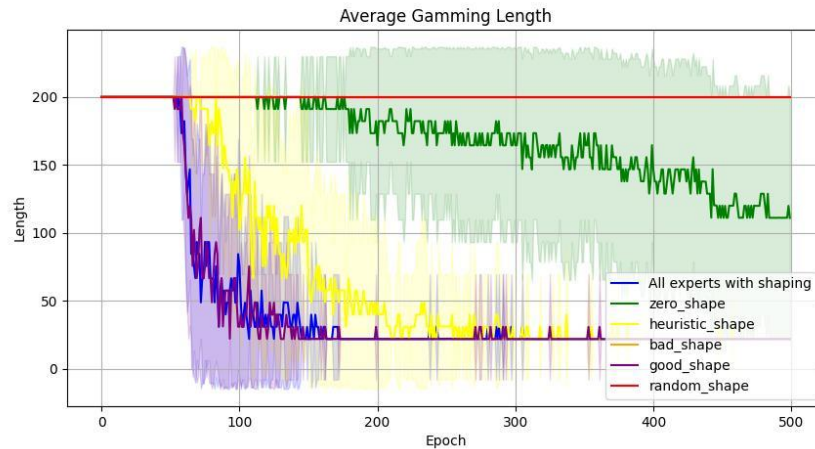
Epsilon: $\max(0.98^t, 0.01)$ / $\max(0.98^t, 0.01)$ / $\max(0.98^t, 0.01)$

Episode: 2000 / 500 / 2000

Independent set: 20 / 20 / 20

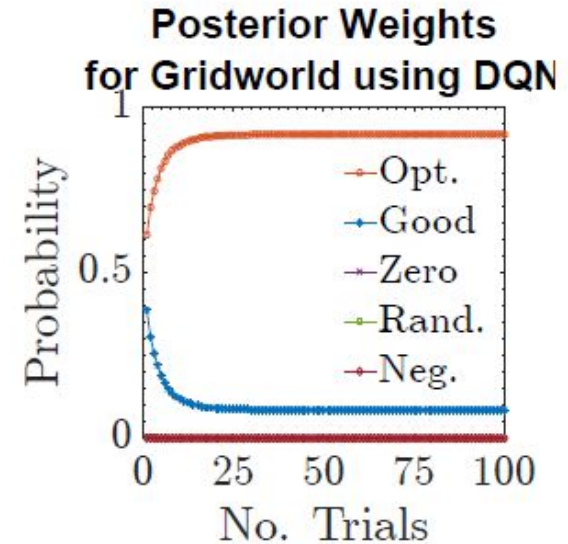
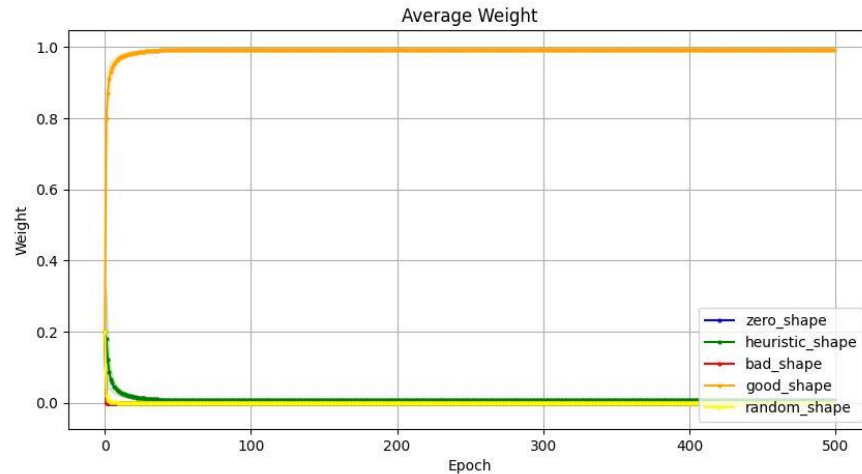
Experimental Results

DQN — Gridworld



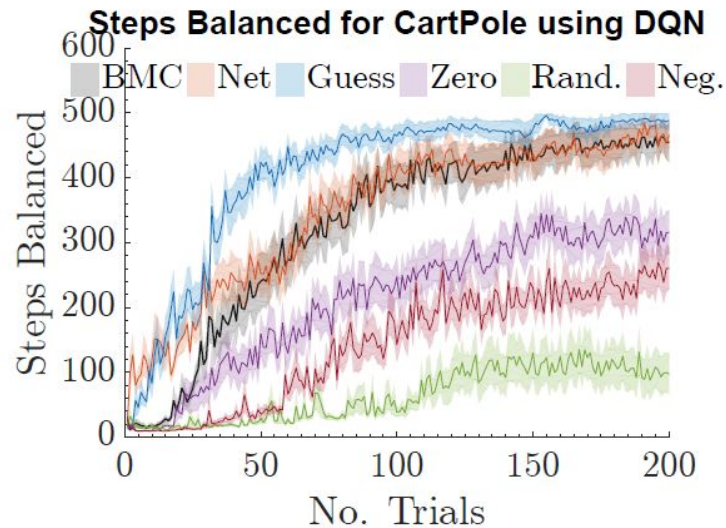
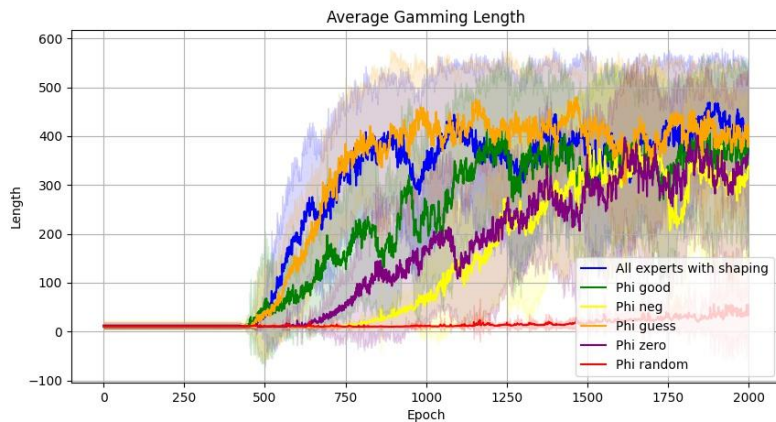
Experimental Results

DQN — Gridworld



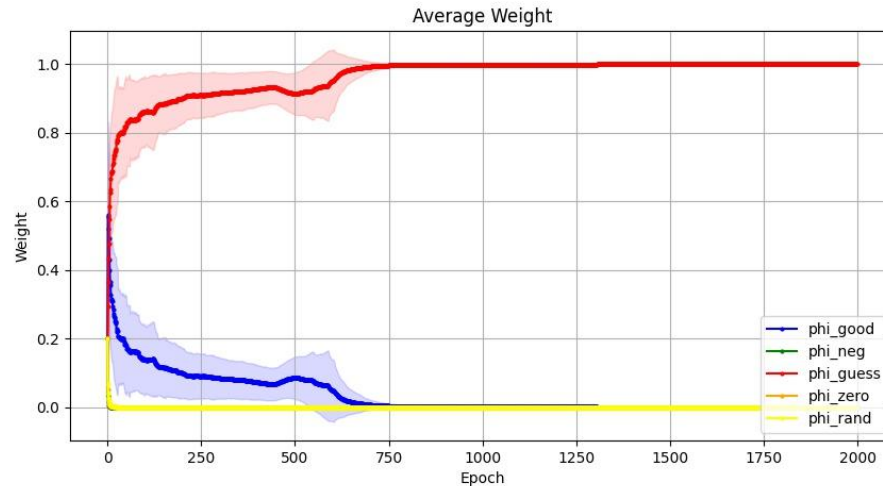
Experimental Results

DQN — CartPole

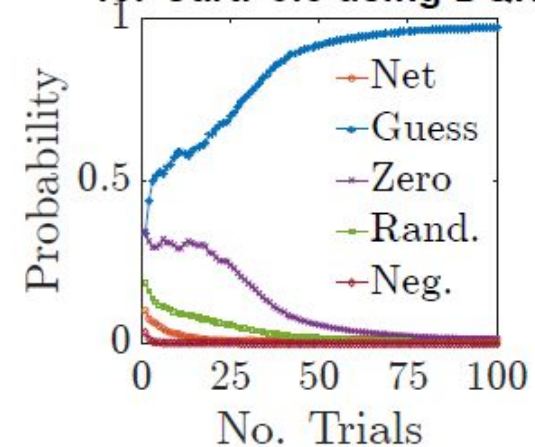


Experimental Results

DQN — CartPole

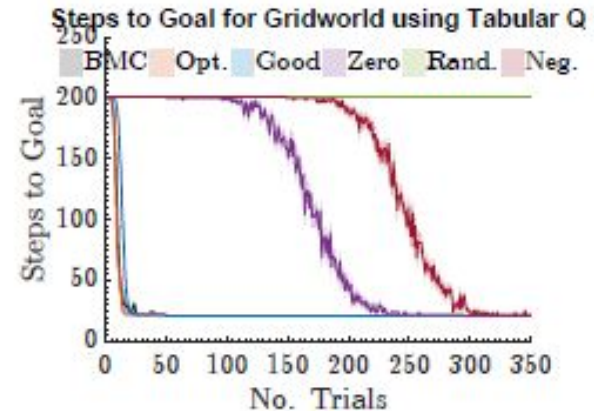
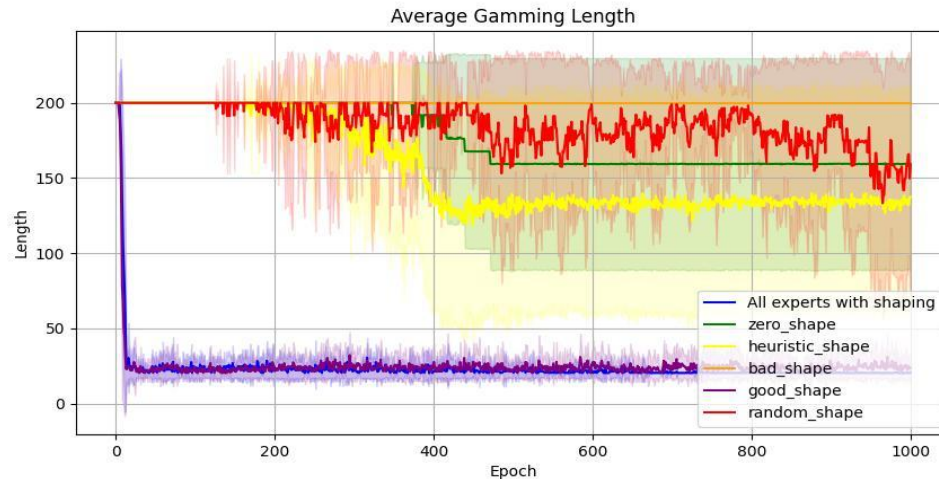


Posterior Weights for CartPole using DQN



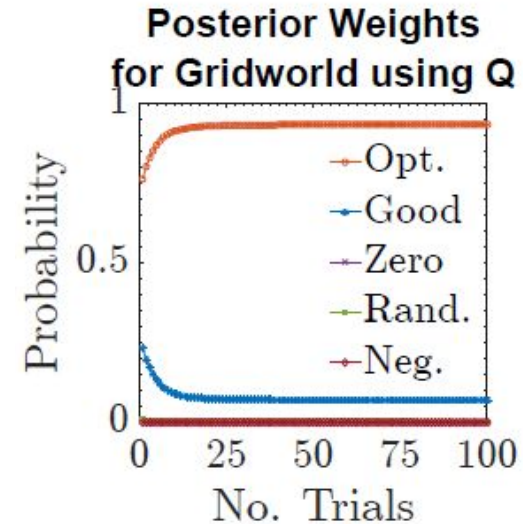
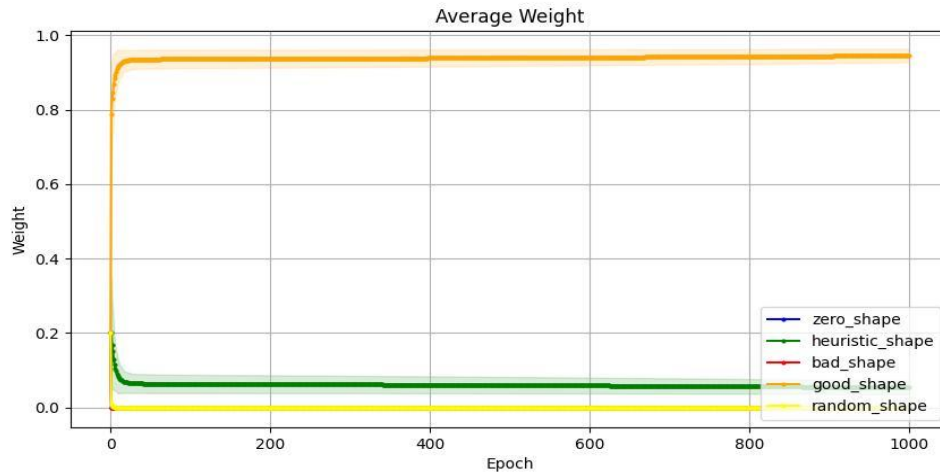
Experimental Results

Q-learning — Gridworld without min_eps



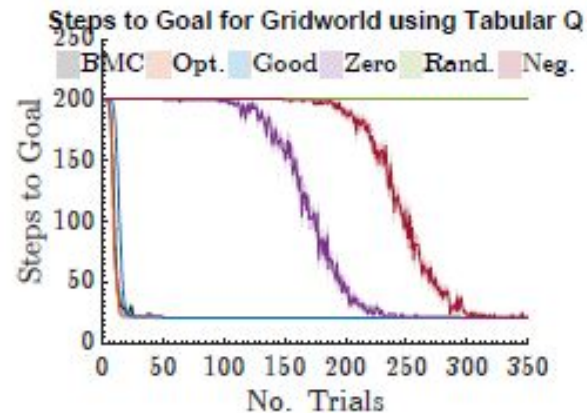
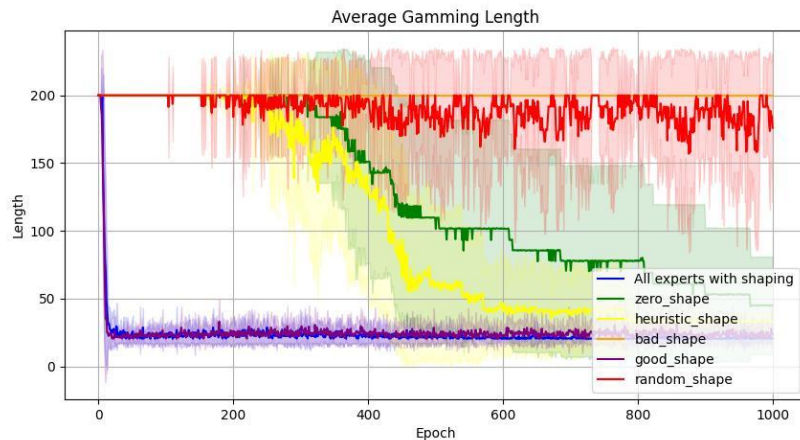
Experimental Results

Q-learning — Gridworld without min_eps



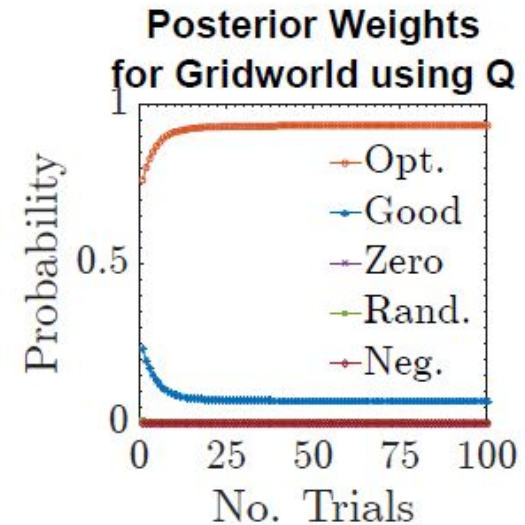
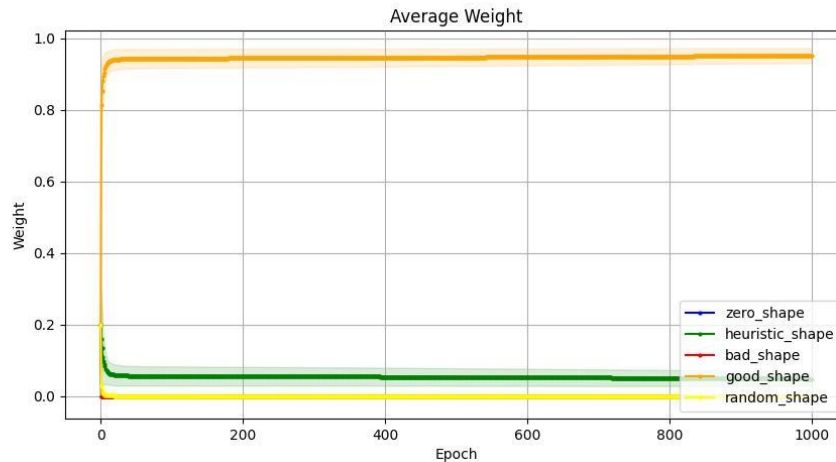
Experimental Results

Q-learning — Gridworld with min_eps



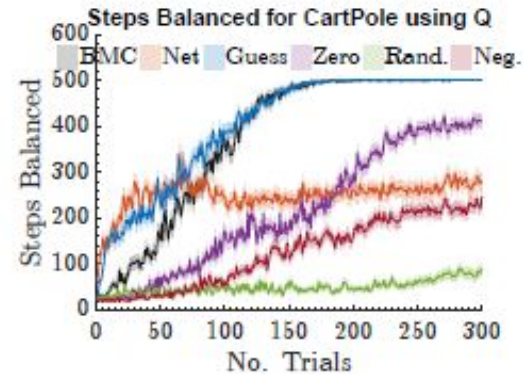
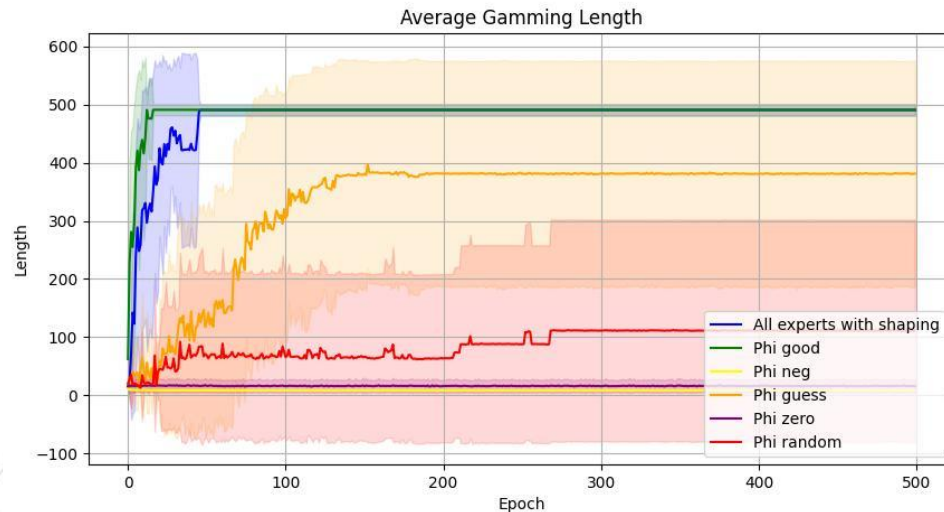
Experimental Results

Q-learning — Gridworld with min_eps



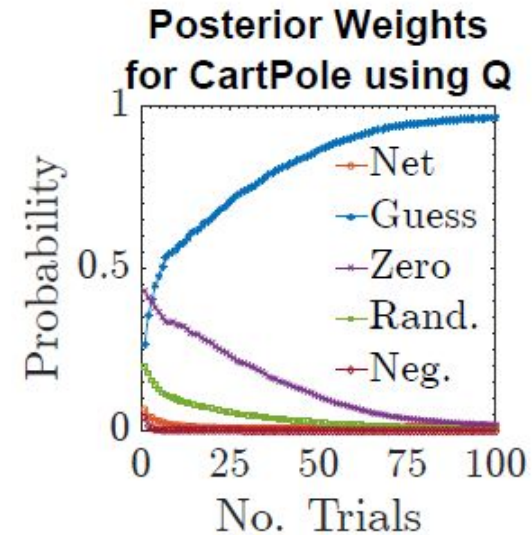
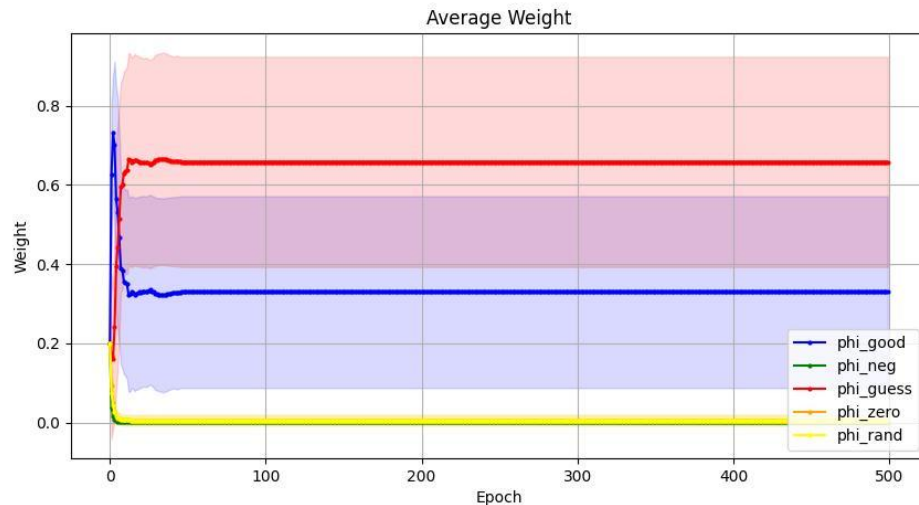
Experimental Results

Q-learning — CartPole with early termination



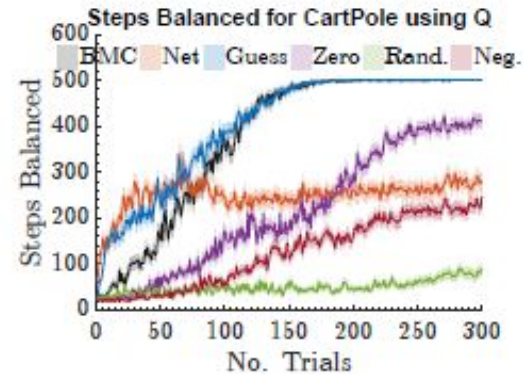
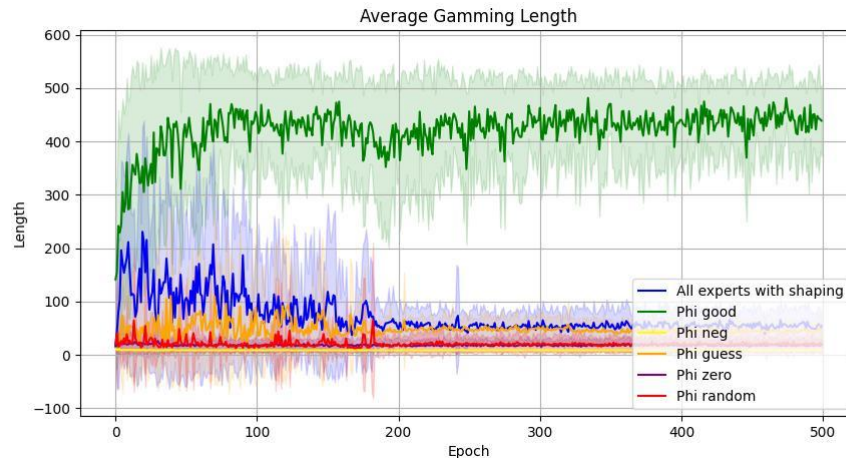
Experimental Results

Q-learning — CartPole with early termination



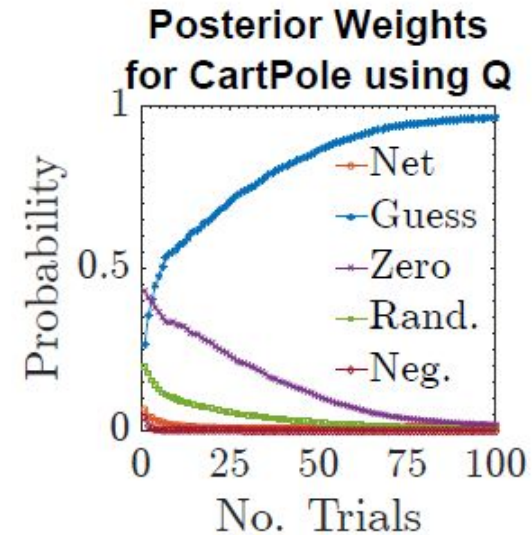
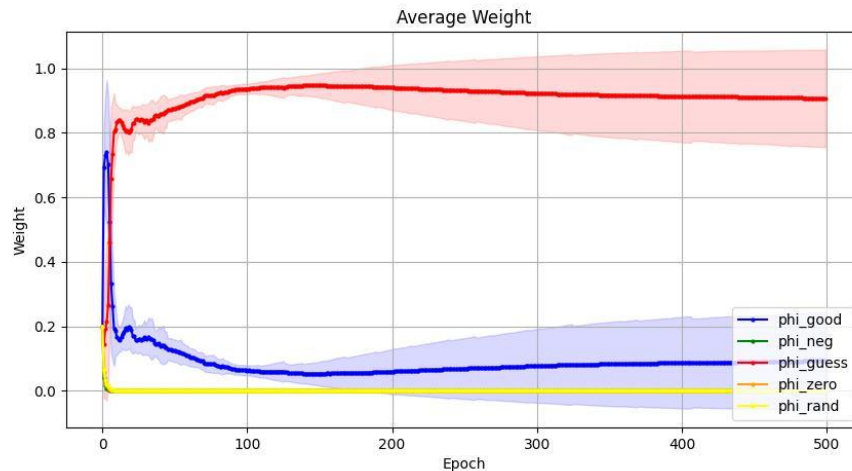
Experimental Results

Q-learning — CartPole without early termination



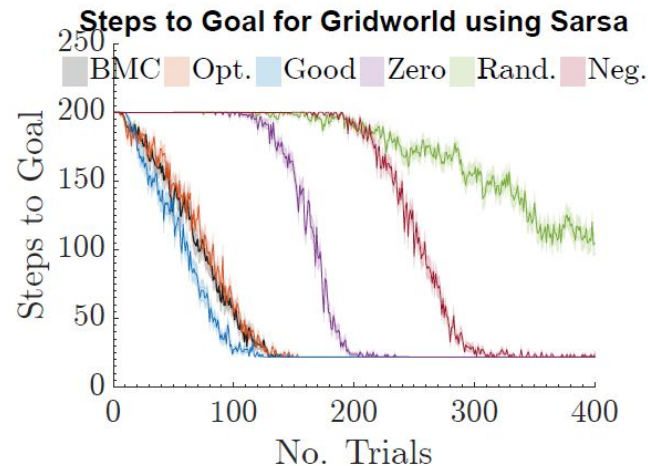
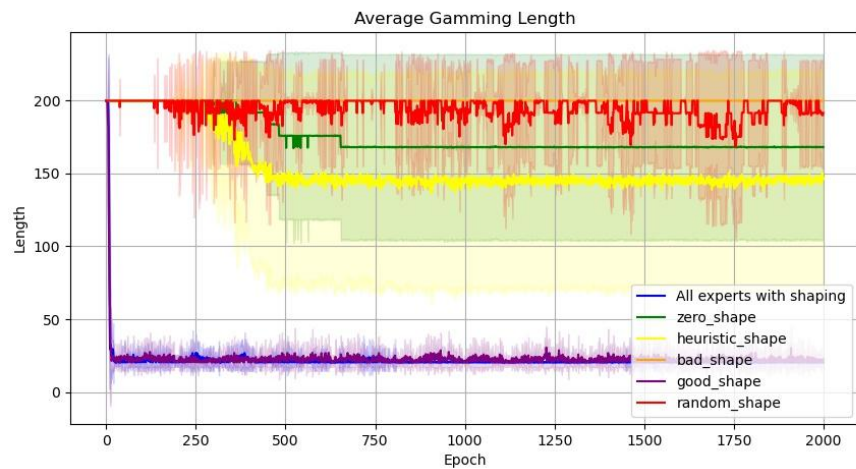
Experimental Results

Q-learning — CartPole without early termination



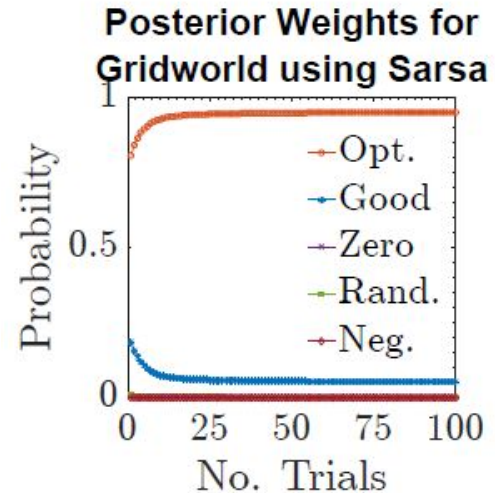
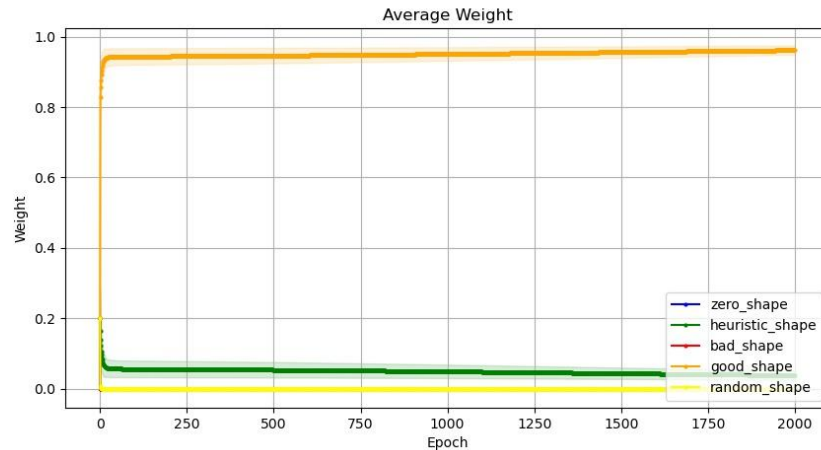
Experimental Results

SARSA — Gridworld without min_eps



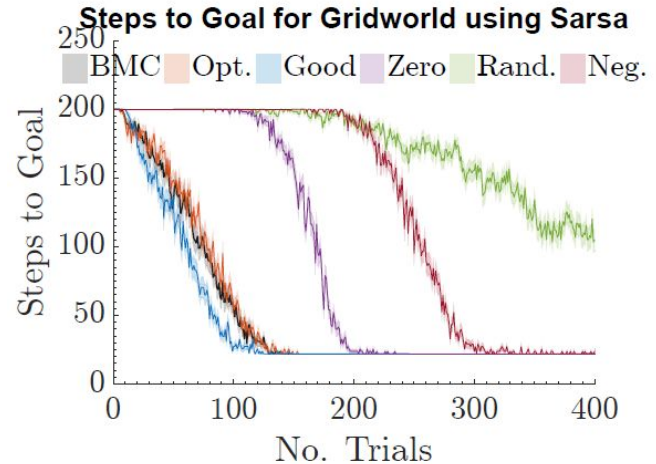
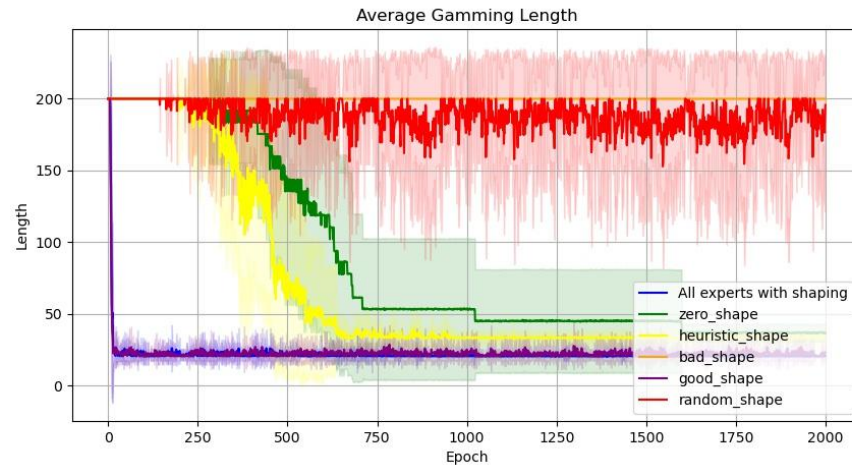
Experimental Results

SARSA — Gridworld without min_eps



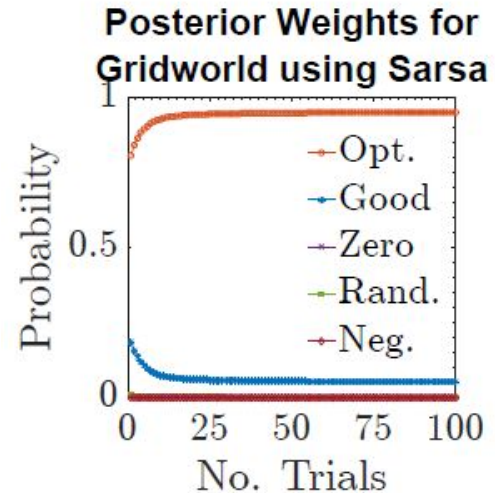
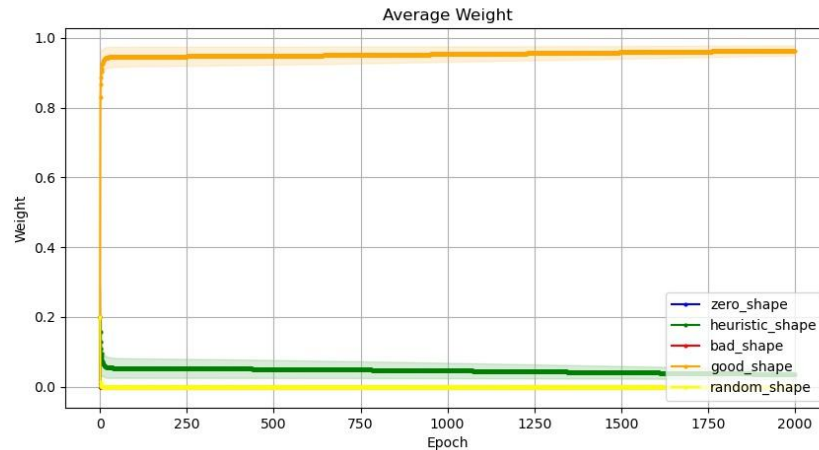
Experimental Results

SARSA — Gridworld with min_eps



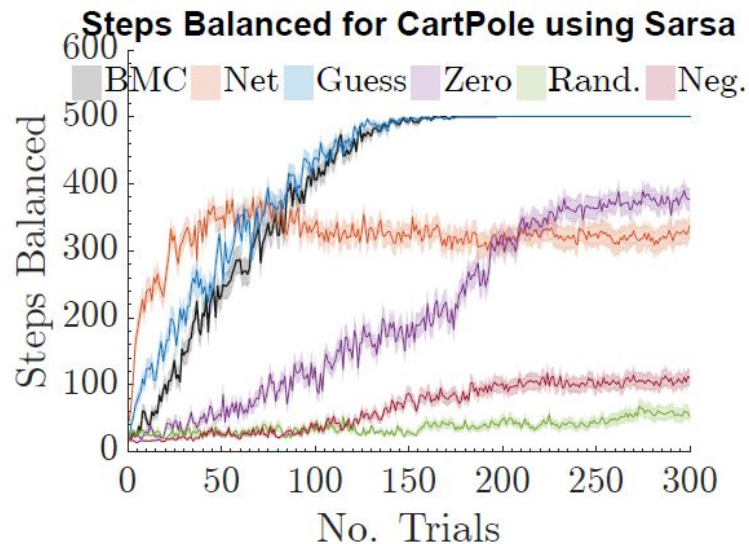
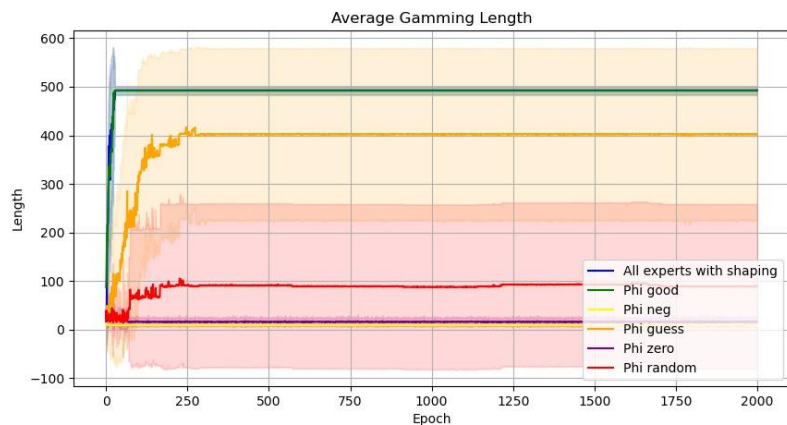
Experimental Results

SARSA — Gridworld with min_eps



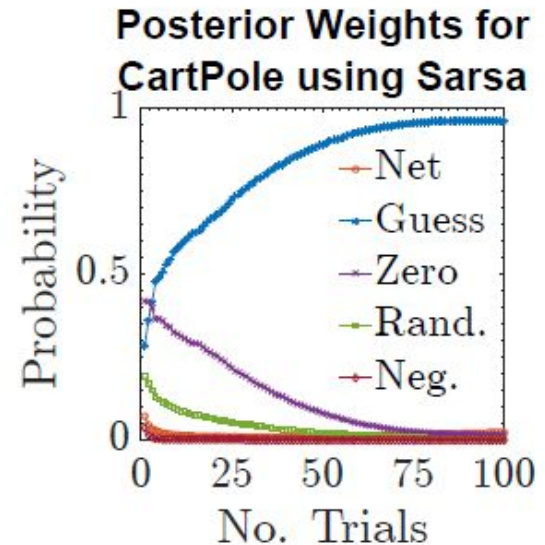
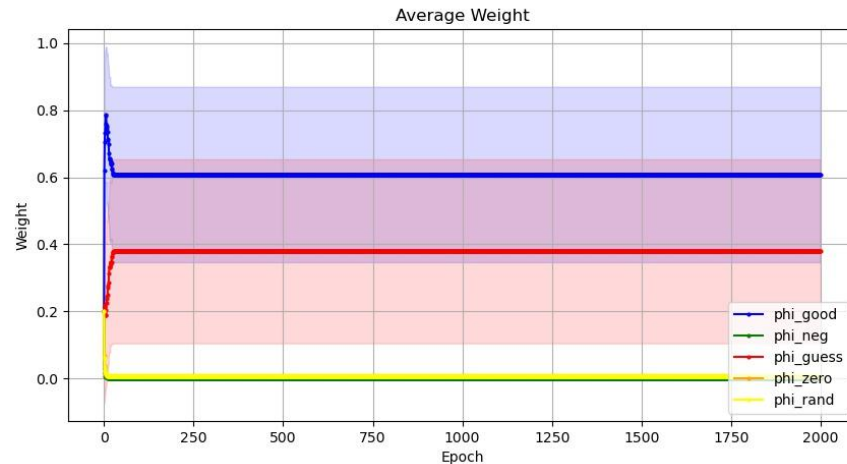
Experimental Results

SARSA — CartPole with early termination



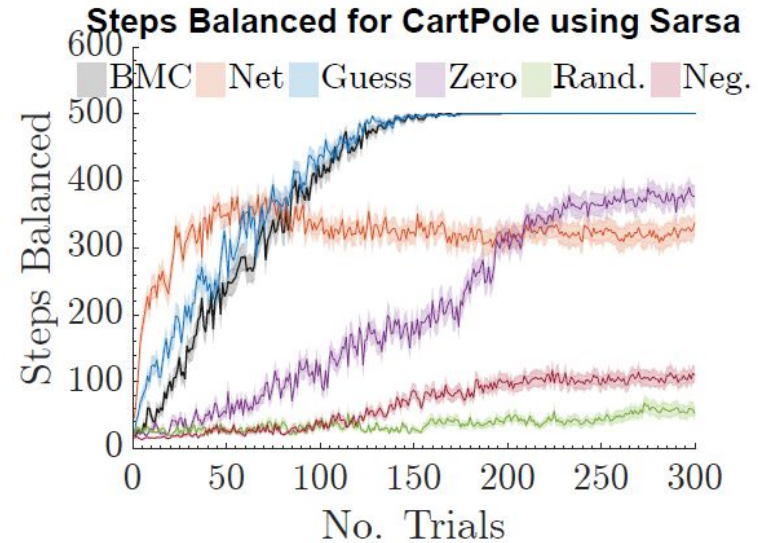
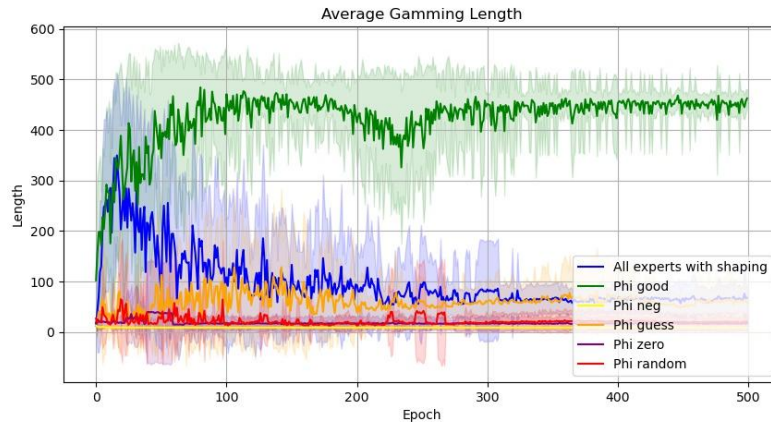
Experimental Results

SARSA — CartPole with early termination



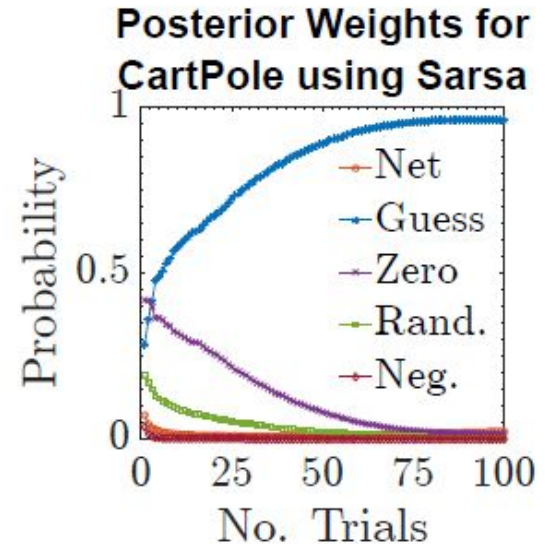
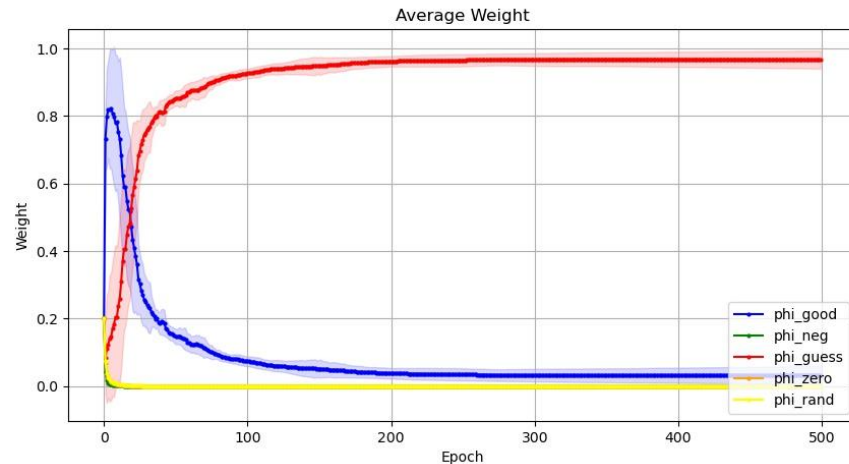
Experimental Results

SARSA — CartPole without early termination



Experimental Results

SARSA — CartPole without early termination



A decorative network diagram in the top-left corner, featuring a complex web of interconnected nodes and lines. The nodes are represented by small circles, some of which are larger and have concentric circles, suggesting a hierarchical or multi-layered structure. The lines are thin and gray, connecting the nodes in a non-linear fashion.

Conclusion

Conclusion

- ◎ Bayesian Reward Shaping isn't always giving us a useful combination of experts for better speeding up the learning process.
- ◎ Early termination might influence the weights provided by Bayesian Reward Shaping.
- ◎ The lower bound of epsilon (for exploration) is important in both tabular methods.



Thanks!

Any questions?