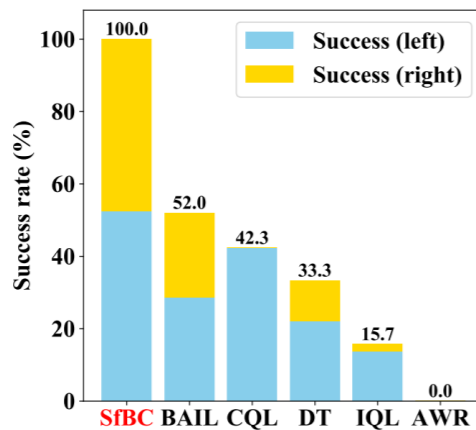
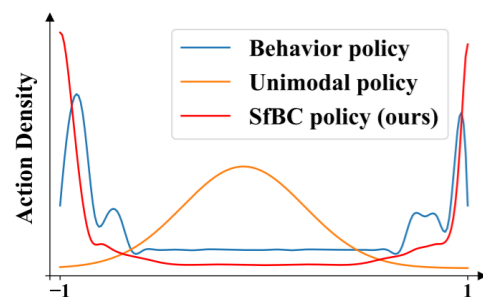


Motivation

- Traditional weighted regression methods generally use Gaussian policies which **lack distributional expressivity**.
- The behavior dataset are usually highly diverse, and the optimal decision space may be **multimodal**.
- Limited expressivity may lead to the OOD problem during dynamic programming.



- Diffusion models** are powerful generative models, which may potentially be helpful to modeling a heterogeneous behavior dataset.

Challenges

- Diffusion models is an implicit generative model, which means that calculation of log probability is not tractable.
- Weighted regression method cannot be directly applied.

$$\arg \max_{\theta} \mathbb{E}_{(s,a) \sim \mathcal{D}^{\mu}} \left[\frac{1}{Z(s)} \log \pi_{\theta}(a|s) \exp(\alpha Q_{\phi}(s,a)) \right]$$

Difficult to analytically calculate

Method

Constrained policy search:

$$\arg \max_{\pi} \int_{\mathcal{S}} \rho_{\mu}(s) \int_{\mathcal{A}} \pi(a|s) Q_{\phi}(s,a) da ds - \frac{1}{\alpha} \int_{\mathcal{S}} \rho_{\mu}(s) D_{\text{KL}}(\pi(\cdot|s) || \mu(\cdot|s)) ds.$$

Optimal solution:

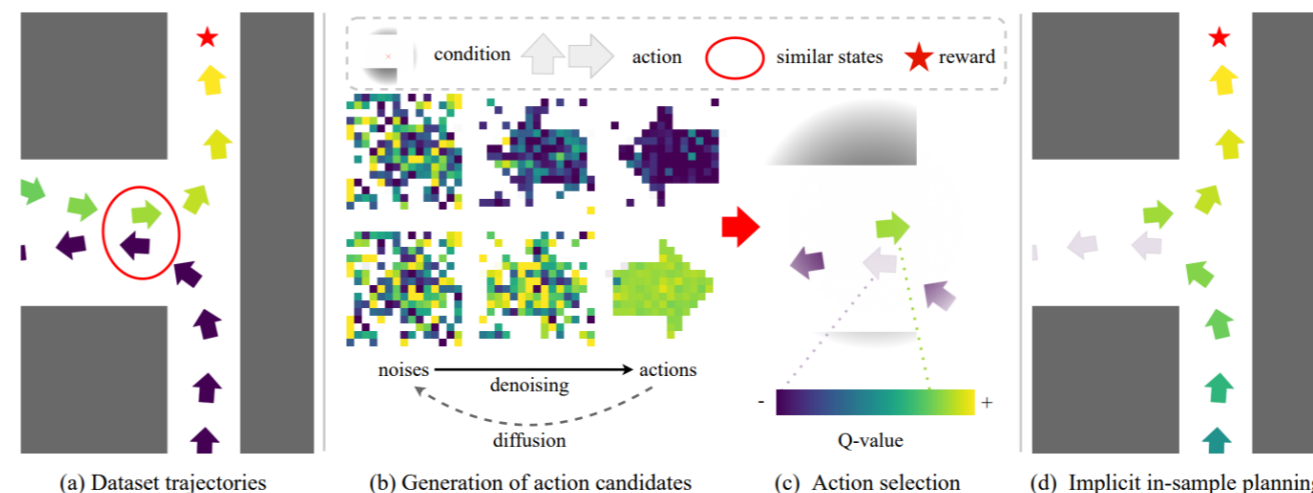
$$\pi^{*}(a|s) = \frac{1}{Z(s)} \mu(a|s) \exp(\alpha Q_{\phi}(s,a))$$

Policy decoupling:

$$\pi(a|s) \propto \mu_{\theta}(a|s) \exp(\alpha Q_{\phi}(s,a))$$

Diffusion modeling:

$$\theta = \arg \min_{\theta} \mathbb{E}_{(s,a) \sim D^{\mu}, \epsilon, t} [\|\sigma_t \mathbf{s} \theta (\alpha_t \mathbf{a} + \sigma_t \epsilon, s, t) + \epsilon\|_2^2]$$



D4RL Experiments

Algorithm	MuJoCo Locomotion	Antmaze	Maze2d	Kitchen
IQL	76.9	63.0	50.0	53.3
BAIL	71.6	46.7	-	-
DT	74.7	18.7	-	-
Diffuser	75.3	-	119.5	-
SfBC (ours)	75.6	74.2	74.0	57.1

Env Decision space

