

Generative Flow Networks

Connections with Reinforcement Learning

Tristan Deleu

Entalpic

Probabilistic inference

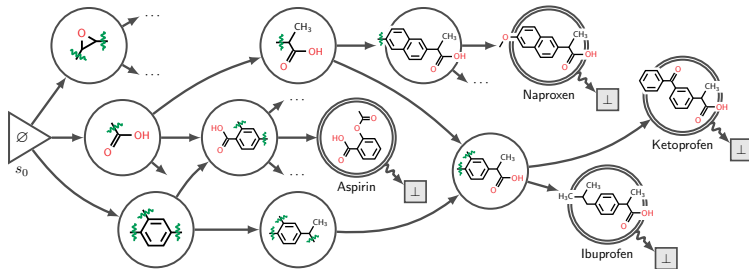
- **Objective:** Sampling from the target distribution

$$P^*(x) = \frac{\exp(-\mathcal{E}(x))}{Z}$$

- States are organized as a **pointed DAG** \mathcal{G}
- **Forward transition probability:**
Distribution over the children of s

$$P_F(\cdot \mid s)$$

- **Sampling** from P_F^\top : follow P_F at random until we reach \perp .



Terminating state distribution

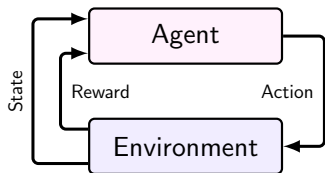
$$P_F^\top(x) \triangleq \sum_{\tau: x \rightarrow \perp \in \tau} P_F(\tau) = \sum_{\tau: s_0 \rightsquigarrow x} \left[\prod_{t=0}^{T_\tau} P_F(s_{t+1} \mid s_t) \right]$$

Outline

- Introduction to Reinforcement Learning
- Connection between GFlowNets and Maximum Entropy RL
- KL-regularized RL
- Intermediate rewards

Reinforcement learning

- **Idea:** An **agent** learns to achieve a goal by interacting with an **environment**.
- The agent takes actions in the environment, and receives rewards. The goal is to **maximize the sum of reward**.



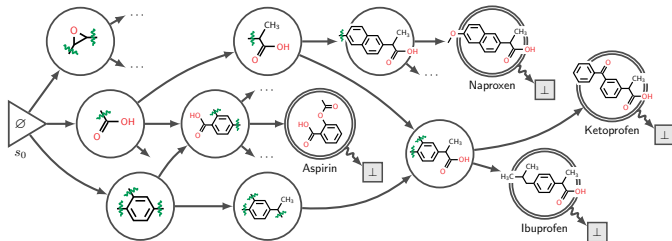
Markov Decision Process

- The environment is represented as a **Markov Decision Process (MDP)**
 - A **state** space \mathcal{S}
 - An **action** space \mathcal{A}
 - A **transition** function $\mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$
 - A reward function

$$r(s_t, s_{t+1})$$

- The agent takes actions following a **policy**

$$\pi(s_{t+1} \mid s_t)$$



Objective

$$\pi_{\text{RL}}^* = \arg \max_{\pi} \mathbb{E}_{\tau} \left[\sum_{t=0}^T r(s_t, s_{t+1}) \right]$$

Value-based methods

- How to find the **optimal policy** π_{RL}^* ?

Optimal policy

$$\pi_{\text{RL}}^*(s' | s) = \begin{cases} 1 & \text{if } s' = \arg \max_{s'' \in \text{Ch}_{\mathcal{G}}(s)} Q^*(s, s'') \\ 0 & \text{otherwise} \end{cases}$$

where $Q^*(s, s')$ satisfies the **Bellman optimality equations**

Bellman optimality equations

$$\begin{aligned} Q^*(s, s') &= r(s, s') + V^*(s') \\ V^*(s) &= \max_{s'' \in \text{Ch}_{\mathcal{G}}(s)} Q^*(s, s'') \end{aligned}$$

- **Idea:** Find the optimal $Q^*(s, s')$, and derive the optimal policy from it.

Inference as control

- Treat probabilistic inference as a **control problem** in a Markov Decision Process

- \mathcal{G} is the structure of the MDP
- The reward function defined such that

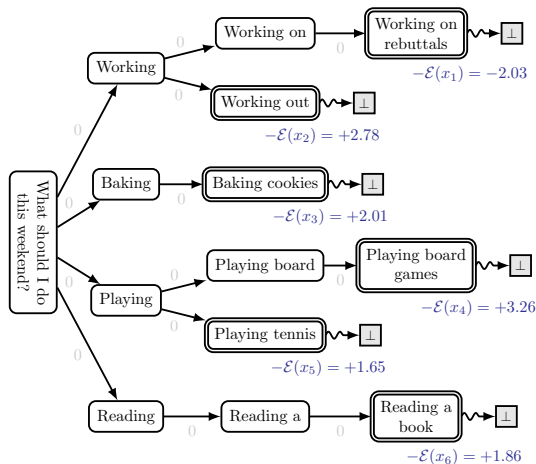
$$\sum_{t=0}^T r(s_t, s_{t+1}) = -\mathcal{E}(s_T)$$

e.g., sparse reward received upon termination.

- Find the **optimal policy**

$$\pi_{\text{RL}}^* = \arg \max_{\pi} \mathbb{E}_{\tau} \left[\sum_{t=0}^T r(s_t, s_{t+1}) \right]$$

- This finds a **mode** of the target distribution
 $P^*(x) \propto \exp(-\mathcal{E}(x))$



Maximum-entropy reinforcement learning

- **Idea:** Add a **regularizer** to encourage diversity

Maximum-entropy reinforcement learning

$$\pi_{\text{MaxEnt}}^* = \arg \max_{\pi} \mathbb{E}_{\tau} \left[\sum_{t=0}^T r(s_t, s_{t+1}) + \mathcal{H}(\pi(\cdot | s_t)) \right]$$

where

$$\mathcal{H}(\pi(\cdot | s)) = - \sum_{s' \in \text{Ch}_{\mathcal{G}}(s)} \pi(s' | s) \log \pi(s' | s)$$

is the **entropy** of the policy π .

- Higher entropy \Rightarrow more *uniform* policy.

RL

MaxEnt RL

Bellman equations

$$Q^*(s, s') = r(s, s') + V^*(s')$$

$$V^*(s) = \max_{s'' \in \text{Ch}_{\mathcal{G}}(s)} Q^*(s, s'')$$

$$Q_{\text{soft}}^*(s, s') = r(s, s') + V_{\text{soft}}^*(s')$$

$$V_{\text{soft}}^*(s) = \log \sum_{s'' \in \text{Ch}_{\mathcal{G}}(s)} \exp(Q_{\text{soft}}^*(s, s''))$$

Optimal policy

$$\pi_{\text{RL}}^*(s' | s) = \arg \max_{s' \in \text{Ch}_{\mathcal{G}}(s)} Q^*(s, s')$$

$$\pi_{\text{soft}}^*(s' | s) = \text{softmax}(Q_{\text{soft}}^*(s, s'))$$

Maximum-entropy reinforcement learning

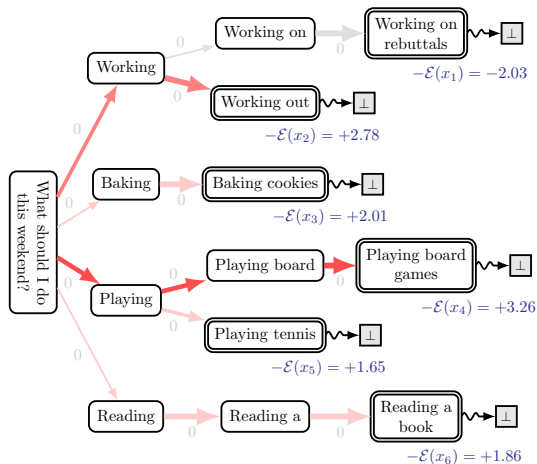
- **Objective:**

$$\pi_{\text{MaxEnt}}^* = \arg \max_{\pi} \mathbb{E}_{\tau} \left[\sum_{t=0}^T r(s_t, s_{t+1}) + \mathcal{H}(\pi(\cdot | s_t)) \right]$$

Theorem

$$\pi_{\text{MaxEnt}}^*(\tau) \propto \exp \left(\sum_{t=0}^T r(s_t, s_{t+1}) \right)$$

- **Recall:** The reward function is defined such that $\sum_{t=0}^T r(s_t, s_{t+1}) = -\mathcal{E}(s_T)$.
- If we know the optimal policy π_{MaxEnt}^* , then we can sample from $P^*(x) \propto \exp(-\mathcal{E}(x))$ by following π_{MaxEnt}^* .



$$\begin{aligned}\pi_{\text{MaxEnt}}^*(\tau) &= \prod_{t=0}^T \pi_{\text{MaxEnt}}^*(s_{t+1} \mid s_t) \\&= \prod_{t=0}^T \exp(r(s_t, s_{t+1}) + V_{\text{soft}}^*(s_{t+1}) - V_{\text{soft}}^*(s_t)) \\&= \exp\left(\sum_{t=0}^T (r(s_t, s_{t+1}) + V_{\text{soft}}^*(s_{t+1}) - V_{\text{soft}}^*(s_t))\right) \\&= \exp\left(\sum_{t=0}^T r(s_t, s_{t+1}) + V_{\text{soft}}^*(\perp) - V_{\text{soft}}^*(s_0)\right) \\&= \frac{1}{\exp(V_{\text{soft}}^*(s_0))} \exp\left(\sum_{t=0}^T r(s_t, s_{t+1})\right) \\&\propto \exp\left(\sum_{t=0}^T r(s_t, s_{t+1})\right)\end{aligned}$$

Theorem

$$\pi_{\text{MaxEnt}}^*(\tau) \propto \exp\left(\sum_{t=0}^T r(s_t, s_{t+1})\right)$$

Soft Bellman optimality equations

$$Q_{\text{soft}}^*(s, s') = r(s, s') + V_{\text{soft}}^*(s')$$

$$V_{\text{soft}}^*(s) = \log \sum_{s'' \in \text{Chg}(s)} \exp(Q_{\text{soft}}^*(s, s''))$$

Soft optimal policy

$$\pi_{\text{MaxEnt}}^*(s' \mid s) = \exp(Q_{\text{soft}}^*(s, s') - V_{\text{soft}}^*(s))$$

Biased sampling

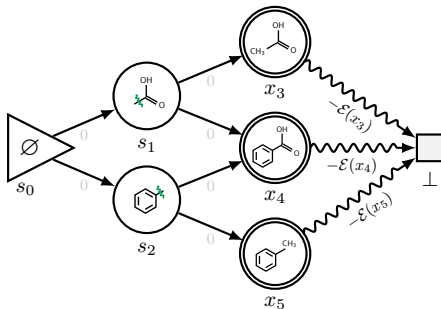
- If $\sum_t r(s_t, s_{t+1}) = -\mathcal{E}(s_T)$, then this is *almost* matching the target distribution $P^*(x) \propto \exp(-\mathcal{E}(x))$.

$$\pi_{\text{MaxEnt}}^*(\tau) \propto \exp\left(\sum_{t=0}^T r(s_t, s_{t+1})\right)$$

- We are interested in a distribution over **outcomes**, not a distribution over **how** these objects are constructed.
- Terminating state distribution:**

$$\pi_{\text{MaxEnt}}^{*\top}(x) = \sum_{\tau: x \rightarrow \perp \in \tau} \pi_{\text{MaxEnt}}^*(\tau) \propto n(x) \exp(-\mathcal{E}(x))$$

- Sampling is **biased** towards terminating states with multiple paths leading to them.



x	τ	$\pi_{\text{MaxEnt}}^*(\tau)$	$\pi_{\text{MaxEnt}}^{*\top}(x)$
x_3	$s_0 \rightarrow s_1 \rightarrow x_3$	$\exp(-\mathcal{E}(x_3))/Z'$	$\exp(-\mathcal{E}(x_3))/Z'$
x_4	$s_0 \rightarrow s_1 \rightarrow x_4$	$\exp(-\mathcal{E}(x_4))/Z'$	$2 \exp(-\mathcal{E}(x_4))/Z'$
	$s_0 \rightarrow s_2 \rightarrow x_4$	$\exp(-\mathcal{E}(x_4))/Z'$	
x_5	$s_0 \rightarrow s_2 \rightarrow x_5$	$\exp(-\mathcal{E}(x_5))/Z'$	$\exp(-\mathcal{E}(x_5))/Z'$

Reward correction

Theorem

Let $P_B(s \mid s')$ be an arbitrary **backward policy** (i.e., a distribution over the parents of $s' \neq s_0$ in \mathcal{G}). If the reward function satisfies for any trajectory $\tau = (s_0, s_1, \dots, s_T, \perp)$:

$$\sum_{t=0}^T r(s_t, s_{t+1}) = -\mathcal{E}(s_T) + \sum_{t=1}^T \log P_B(s_{t-1} \mid s_t),$$

then the terminating state distribution associated with the optimal MaxEnt RL policy satisfies

$$\pi_{\text{MaxEnt}}^{\star\top}(x) \propto \exp(-\mathcal{E}(x))$$

Proof

$$\begin{aligned}\pi_{\text{MaxEnt}}^{\star\top}(x) &= \sum_{\tau:s_0 \rightsquigarrow x} \prod_{t=0}^{T_\tau} \pi_{\text{MaxEnt}}^{\star}(s_{t+1} \mid s_t) \\&= \sum_{\tau:s_0 \rightsquigarrow x} \exp \left[\sum_{t=0}^{T_\tau} (r(s_t, s_{t+1}) + V_{\text{soft}}^{\star}(s_{t+1}) - V_{\text{soft}}^{\star}(s_t)) \right] \\&= \sum_{\tau:s_0 \rightsquigarrow x} \exp \left[\sum_{t=0}^{T_\tau} r(s_t, s_{t+1}) + V_{\text{soft}}^{\star}(\perp) - V_{\text{soft}}^{\star}(s_0) \right] \\&= \sum_{\tau:s_0 \rightsquigarrow x} \exp \left[-\mathcal{E}(x) + \sum_{t=1}^{T_\tau} \log P_B(s_{t-1} \mid s_t) - V_{\text{soft}}^{\star}(s_0) \right] \\&= \exp(-\mathcal{E}(x) - V_{\text{soft}}^{\star}(s_0)) \underbrace{\sum_{\tau:s_0 \rightsquigarrow x} \prod_{t=1}^{T_\tau} P_B(s_{t-1} \mid s_t)}_{=1} \\&= \frac{1}{\exp(V_{\text{soft}}^{\star}(s_0))} \exp(-\mathcal{E}(x))\end{aligned}$$

Corrected reward

$$\sum_{t=0}^T r(s_t, s_{t+1}) = -\mathcal{E}(s_T) + \sum_{t=1}^T \log P_B(s_{t-1} \mid s_t)$$

Soft Bellman optimality equations

$$Q_{\text{soft}}^{\star}(s, s') = r(s, s') + V_{\text{soft}}^{\star}(s')$$

$$V_{\text{soft}}^{\star}(s) = \log \sum_{s'' \in \text{Ch}_G(s)} \exp(Q_{\text{soft}}^{\star}(s, s''))$$

Soft optimal policy

$$\pi_{\text{MaxEnt}}^{\star}(s' \mid s) = \exp(Q_{\text{soft}}^{\star}(s, s') - V_{\text{soft}}^{\star}(s))$$

Examples

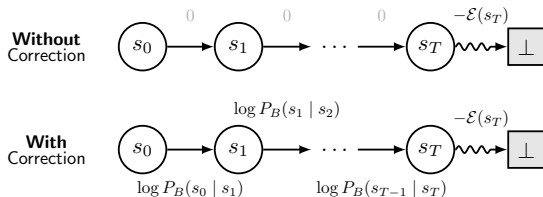
- Corrected reward:

$$\sum_{t=0}^T r(s_t, s_{t+1}) = \mathcal{E}(s_T) + \sum_{t=1}^T \log P_B(s_{t-1} | s_t)$$

- In the case of a sparse reward, we can **add the correction along the way**.
- If the MDP is a **tree**, then $P_B(s | s') = 1$: we recover the case with uncorrected rewards.
- If the backward policy is

$$P_B(s_t | s_{t+1}) = n(s_t)/n(s_{t+1})$$

then we get a policy **maximizing the entropy** over trajectories.



$$r(s_t, s_{t+1}) = \log P_B(s_t | s_{t+1})$$
$$r(s_T, \perp) = -\mathcal{E}(s_T)$$

Path Consistency Learning

- **Recall:** Soft Bellman optimality equations & optimal policy

$$Q_{\text{soft}}^*(s, s') = r(s, s') + V_{\text{soft}}^*(s')$$

$$\pi_{\text{MaxEnt}}^*(s' | s) = \exp(Q_{\text{soft}}^*(s, s') - V_{\text{soft}}^*(s))$$

- Unrolling $\log \pi_{\text{MaxEnt}}^*(s_{t+1} | s_t)$ along a trajectory, we get

Path Consistency Learning

$$\sum_{t=0}^T \log \pi_{\text{MaxEnt}}^*(s_{t+1} | s_t) = \sum_{t=0}^T r(s_t, s_{t+1}) - V_{\text{soft}}^*(s_0)$$

Trajectory Balance

$$\log \frac{Z \prod_{t=0}^T P_F(s_{t+1} | s_t)}{\exp(-\mathcal{E}(s_T)) \prod_{t=1}^T P_B(s_{t-1} | s_t)}$$

- PCL is equivalent to TB, with the reward correction $\sum_{t=0}^T r(s_t, s_{t+1}) = -\mathcal{E}(s_T) + \sum_{t=1}^T \log P_B(s_{t-1} | s_t)$
- There exist other equivalences with MaxEnt RL algorithms (e.g., FL-DB & Soft Q-Learning).

KL-regularized Reinforcement Learning

- We can use **other types of regularizers**. For example if we want the policy to not deviate too much from a **base policy** $\mu(\cdot \mid s)$.

KL-regularized Reinforcement Learning

$$\pi_{\text{RelEnt}}^* = \arg \max_{\pi} \mathbb{E}_{\tau} \left[\sum_{t=0}^T r(s_t, s_{t+1}) - \text{KL}(\pi(\cdot \mid s_t) \parallel \mu(\cdot \mid s_t)) \right]$$

where

$$\text{KL}(\pi(\cdot \mid s) \parallel \mu(\cdot \mid s)) = \sum_{s' \in \text{Ch}_{\mathcal{G}}(s)} \pi(s' \mid s) \log \frac{\pi(s' \mid s)}{\mu(s' \mid s)}$$

- The policy $\mu(\cdot \mid s)$ may be a **prior distribution** (e.g., a distribution learned from data)
- This type of objective is very popular recently with **RL-finetuning of LLMs**.

Terminating state distribution

$$\pi_{\text{RelEnt}}^* = \arg \max_{\pi} \mathbb{E}_{\tau} \left[\sum_{t=0}^T r(s_t, s_{t+1}) - \text{KL}(\pi(\cdot \mid s_t) \parallel \mu(\cdot \mid s_t)) \right]$$

Theorem

If the reward function satisfies for any trajectory $\tau = (s_0, s_1, \dots, s_T, \perp)$:

$$\sum_{t=0}^T r(s_t, s_{t+1}) = -\mathcal{E}(s_T),$$

then the terminating state distribution associated with the optimal KL-regularized RL policy satisfies

$$\pi_{\text{RelEnt}}^{\star \top}(x) \propto \mu^{\top}(x) \exp(-\mathcal{E}(x))$$

Relative Trajectory Balance

- Modification of the **Trajectory Balance** to account for a prior distribution

Relative Trajectory Balance

$$Z \prod_{t=0}^T P_F(s_{t+1} \mid s_t) = \exp(-\mathcal{E}(s_T)) \prod_{t=0}^T \pi_{\text{prior}}(s_{t+1} \mid s_t)$$

- If RTB is satisfied for all trajectories, then $P_F^\top(x) \propto \pi_{\text{prior}}^\top(x) \exp(-\mathcal{E}(x))$.
- This is **equivalent** to a KL-regularized RL algorithm called **Trust-PCL** (without any reward correction needed).

Energy difference as intermediate reward

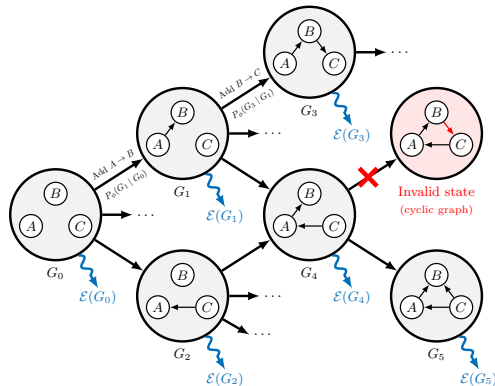
- In cases where **all the states are terminating**, then we can use the difference in energies as an intermediate reward

$$r(s_t, s_{t+1}) = \mathcal{E}(s_t) - \mathcal{E}(s_{t+1})$$

this guarantees that

$$\sum_{t=0}^T r(s_t, s_{t+1}) = \underbrace{\mathcal{E}(s_0)}_{\text{constant}} - \mathcal{E}(s_T)$$

- Examples:** Hypergrid, structure learning.
- Consequence:** FL-DB does not depend on a flow (value) function anymore. So-called “**modified**” DB.



Diffusion model

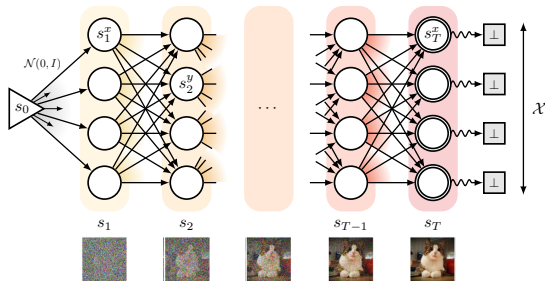
- Noising process $x_T \rightarrow x_{T-1} \rightarrow \dots \rightarrow x_1$

$$q(x_{t-1} | x_t) = \mathcal{N}(x_{t-1}; \sqrt{\alpha_t}x_t, (1 - \alpha_t)I)$$

- Objective:** Find a model $p_\theta(x_{t+1} | x_t)$ that denoises x_t to x_{t+1} .

$$p_\theta(x_{t+1} | x_t) = \mathcal{N}(x_{t+1}; \mu_\theta(x_t, t), \sigma_t^2 I)$$

where we start with $p(x_1) = \mathcal{N}(x_1; 0, I)$.



	Diffusion model	GFlowNet	First step
Denoising	$p_\theta(x_{t+1} x_t)$	$P_F(x_{t+1} x_t) = \mathcal{N}(\mu_\theta(x_t, t), \sigma_t^2 I)$	$\mathcal{N}(0, I)$
Noising	$q(x_{t-1} x_t)$	$P_B(x_{t-1} x_t) = \mathcal{N}(\sqrt{\alpha_t}x_t, (1 - \alpha_t)I)$	δ_{s_0}

Steering diffusion models

- **Diffusion sampler:** Steering the generation of images with a diffusion process based on an energy $\mathcal{E}(x)$ i.e., sampling from the distribution $\propto \exp(-\mathcal{E}(x))$.

- Parametrization of the denoiser / forward policy

$$P_F(x_{t+1} \mid x_t) = \mathcal{N}(\mu_\theta(x_t, t), \sigma_t^2 I)$$

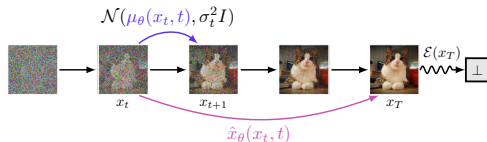
$$\mu_\theta(x_t, t) = a_t x_t + b_t \hat{x}_\theta(x_t, t)$$

- Use the x_T -prediction model as a **proxy**, and evaluate the difference in energies on the predictions

$$r(s_t, s_{t+1}) = \mathcal{E}(\hat{x}_\theta(s_t)) - \mathcal{E}(\hat{x}_\theta(s_{t+1}))$$

With careful treatment of edge cases, this guarantees

$$\sum_{t=0}^T r(s_t, s_{t+1}) = -\mathcal{E}(x_T)$$



- **Idea #1:** Use Trajectory Balance
- **Idea #2:** Use the difference in energies
 $r(s_t, s_{t+1}) = \mathcal{E}(x_t) - \mathcal{E}(x_{t+1})$

Conclusion

- There are **strong connections** between GFlowNets and Reinforcement Learning
- All the techniques from RL can be transferred to GFlowNet applications:
 - Replay buffers
 - Target networks
 - Distributional
 - Advanced exploration strategies
- **Takeaways:**
 - **Be careful** of papers claiming the GFlowNets are significantly better than RL
 - Using GFlowNets for autoregressive tasks (e.g., LLMs) is **exactly equivalent** to MaxEnt RL