

## Problem Formulation

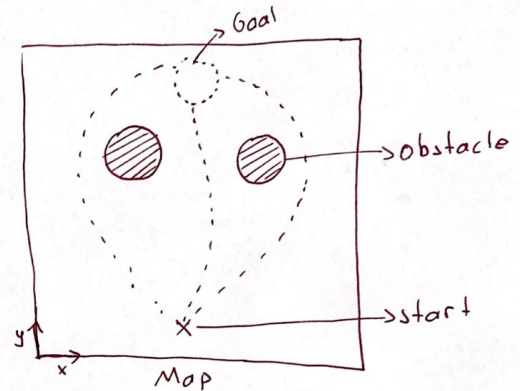
=> For a given environment with obstacles, find a safe path between a starting point and a goal position without any collisions.

# Robot Dynamics

$$X_{k+1} = AX_k + Bu_k + w_k$$

$$x_k \in \mathcal{X} \subset \mathbb{R}^2: [x, y] \text{ position of the robot}$$
$$v_k \in A \subset \mathbb{R}^2 : \text{Input}$$

$w_k \in \mathbb{R}^2$  : Process disturbance with unknown distribution  $P_{true}$



## Markov Decision Problem

 $\langle S, A, r, P \rangle$ 

$S \subset \mathbb{R}^n$ : Bounded cont.  $\Rightarrow s_k = [x_k, o_k^{(ci)}, g]$   
 state space                      Position      Distance to obstacle centroids      goal position  
 $A \subset \mathbb{R}^2$ : Finite action space

$A \subset \mathbb{R}^2$ : Finite action space

$$r: S \times A \times S \rightarrow \mathbb{R}: \text{Reward function}$$

$P: S \times A \rightarrow \Delta_S$ : Transition probability  
 $P(s_{k+1} | s_k, a_k)$

$$\pi: S \rightarrow \Delta_A : \text{Policy, } P(a_k | s_k)$$

Trajectory:  $\mathcal{T} = (s_0, \sigma_0, s_1, \sigma_1, \dots, s_N)$

Discounted Returns:  $R(\pi) = \sum_{i=0}^{N-1} \gamma^i r(s_i, a_i, s_{i+1})$

Value Function:  $V^\pi(s) = \mathbb{E}_{a_t \sim \pi, s_t \sim p} [R(\mathcal{T}^\pi | s_0 = s)]$

Q-Function:  $Q(s, a) = E_{\substack{a_t \sim \pi \\ s_{t+1} \sim P}} [R(s^x | s_0=s, a_0=a)]$

The goal

$$\pi^* = \arg \max_{\pi \in \Pi} V^\pi(s), \forall s \in S$$

## Distributionally Robust Optimization

$h(x, z)$ : objective function

$x \in \mathcal{X}$ : controlled variable

$z \in \mathcal{Z}$ : stochastic variable,

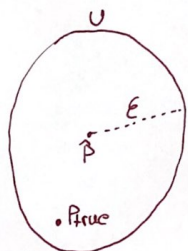
$$\text{DRO: } \inf_{x \in \mathcal{X}} \sup_{P \in \mathcal{U}} \underbrace{E_{z \sim P}[h(x, z)]}_{\text{worst-case scenario}}$$

$\mathcal{U}$ : ambiguity set

$$\mathcal{U} := \{P \mid W_P(\hat{P}, P) \leq \epsilon\}$$

select  $\epsilon$  such that

$$P(P_{\text{true}} \in \mathcal{U}) \geq 1 - \beta$$



## Wasserstein Distance

= Earth movers distance

$$W_P(Q, Q') = \left( \inf_{\pi \in \Pi(Q, Q')} \int_{\mathcal{Z} \times \mathcal{Z}} \|z - z'\|^P \pi(dz, dz') \right)^{1/P}$$

$\Rightarrow$  Min energy required to convert  $Q$  into  $Q'$

$\Pi(Q, Q')$ : set of all joint prob. distributions with marginals  $Q, Q'$

$$\text{Primal: } v_P = \sup_{P \in \mathcal{U}} \int_{\mathcal{Z}} h(z) dP(z) : W_P(P, \hat{P}) \leq \epsilon$$

$\Rightarrow$  The primal problem has a strong dual

$$\text{Dual: } v_D = \inf_{\lambda \geq 0} \left\{ \lambda \epsilon^P - \int_{\mathcal{Z}} \inf_{z \in \mathcal{Z}} [\lambda \|z - z_0\|^P - h(z)] d\hat{P}(z_0) \right\}$$

$$v_P = v_D$$

Discrete Empirical  $\hat{P}_N$

$\Rightarrow N$  samples of  $z$

$$\hat{P}_N \triangleq \frac{1}{N} \sum_{i=1}^N \delta_{z_i}(z)$$

$$v_P = v_D = \min_{\lambda \geq 0} \left\{ \lambda \epsilon^P + \frac{1}{N} \sum_{i=1}^N \sup_{z \in \mathcal{Z}} [h(z) - \lambda \|z - z_i\|^P] \right\}$$

## Distributionally Robust Q-Learning

$$Q\text{-function: } Q(s, a) = E \left[ \sum_{k=0}^{\infty} \gamma^k r(s_k, a_k, s_{k+1}) \mid s_0 = s, a_0 = a \right]$$

### Bellman Operator $\mathcal{T}$

$$\mathcal{T}Q(s, a) = E_{s' \sim p} [r(s, a, s') + \gamma E_{a' \sim \pi} [Q(s', a')]]$$

$\Rightarrow$  Contraction

$$\pi(a|s) = \frac{e^{Q(s, a)}}{\sum_{a' \in A} e^{Q(s, a')}} \quad \sum_{a' \in A} \pi(a'|s) Q(s', a')$$

### Distributionally Robust Bellman Operator $\hat{\mathcal{T}}$

$$\hat{\mathcal{T}}Q(s, a) = \inf_{p \in \mathcal{U}} E_{s' \sim p} \left[ \underbrace{r(s, a, s') + \gamma E_{a' \sim \pi} [Q(s', a')]}_{h(s, a, s') \Rightarrow h(s')} \right]$$

$\Rightarrow$  Worst case expected discounted returns

$h(s') \Rightarrow$  non-linear, non-convex

$\Rightarrow Q(s, a) \rightarrow$  approximated by a neural network



### TD-Learning

$\Rightarrow$  store transition  $(s_k, a_k, s_{k+1}, r_k)$   
 td-error  $\Rightarrow \delta = r_k + \gamma \max_a Q(s_{k+1}, a) - Q(s_k, a_k)$   
 $Q(s_k, a_k) \leftarrow Q(s_k, a_k) + \alpha \delta$  (target  $s$ )

Distributionally Robust target  $= \inf_{p \in \mathcal{U}} E_{s' \sim p} [h(s_k, a_k, s')]$

$\Rightarrow$  Double DQN

$\Rightarrow$  Dueling layer

$\Rightarrow$  Prioritized experience replay



## Solving DRO

### Lipschitz Approximation

$$\star \sup_{P \in \mathcal{U}} E_{s \sim P}[-h(s)]$$

$$\text{Lip}(h) = \sup_{\{s, s'\}} \frac{|h(s) - h(s')|}{\|s - s'\|}$$

$$\Rightarrow \sup_{P \in \mathcal{U}} E_{s \sim P}[-h(s)] \leq E_{s \sim \hat{P}_N}[-h(s)] + \epsilon \text{Lip}(h)$$

$\Rightarrow \text{Lip}(h)$  can be estimated by using

LipSDP in python

$\hookrightarrow$  neural network lipschitz const. estimation

$\Rightarrow \text{Lip}(h)$  can be estimated when the target network is aligned.

### Robust Program Approximation

$$V_P = V_D = \min_{\lambda \geq 0} \left\{ \lambda \epsilon^P + \frac{1}{N} \sum_{i=1}^N \sup_{z \in \mathcal{Z}} [-h(z) - \lambda \|z - z_i\|^P] \right\}$$

for  $K > 0$

$$V_K := \sup_{(z^{ik})_{i,k} \in \mathcal{M}_K} \frac{1}{NK} \sum_{i=1}^N \sum_{k=1}^K -h(z^{ik})$$

$$\mathcal{M}_K := \left\{ (z^{ik})_{i,k} : \frac{1}{NK} \sum_{i=1}^N \sum_{k=1}^K \|z^{ik} - z_i\|^P \leq \epsilon^P, z^{ik} \in \mathcal{Z}, \forall i, k \right\}$$

$$V_K \uparrow \sup_{P \in \mathcal{U}} E_{s \sim P}[-h(s)] \text{ as } K \rightarrow \infty$$

$\Rightarrow$  Local optimums can be found

$\Rightarrow$  Must solve for each transition in memory

$\Rightarrow$  Will be slow