# 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

$X_{k+1} = A X_k + B U_k + W_k$

$X_k \in \mathcal{X} \subset \mathbb{R}^2$ : $[x, y]$ position of the robot

$U_k \in A \subset \mathbb{R}^2$ : Input

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

# Markov Decision Problem
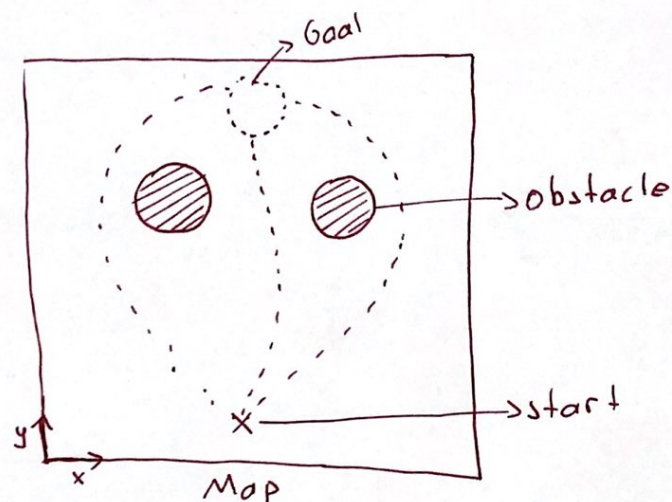
$\langle S, A, r, P \rangle$

$S \subset \mathbb{R}^n$ : Bounded cont. => $s_k = [x_k, o_k^{(i)}, g]$
state space          Position Distance  goal
                              to obstacle Position
                              centroids

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

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

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

$\pi : S \to \Delta_A$ : Policy, $P(a_k | s_k)$



Map

Trajectory: $\mathcal{T}^\pi = (s_0, a_0, s_1, a_1, \ldots, s_N)$

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

Value Function: $V^\pi(s) = \underset{\substack{a_k \sim \pi \\ s_{k+1} \sim P}}{E} [R(\mathcal{T}^\pi | s_0 = s)]$

Q-Function: $Q(s, a) = \underset{\substack{a_k \sim \pi \\ s_{k+1} \sim P}}{E} [R(\mathcal{T}^\pi | s_0 = s, a_0 = a)]$

# The goal

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

# Distributionally Robust Optimization

$h(x,z)$ : objective function

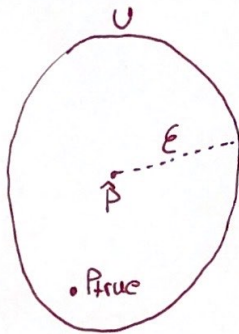$\quad x \in X$ : controlled variable

$\quad z \in Z$ : stochastic variable,

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

$U$ : ambiguity set

$$U := \{P \mid W_p(\hat{P}, P) \leq \varepsilon\}$$

select $\varepsilon$ such that

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

## Wasserstein Distance

= Earth movers distance

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

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

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

Primal: $v_P = \sup_{P \in U} \int_Z h(z) \, dP(z) \; : \; W_p(P, \hat{P}) \leq \varepsilon$

$\Rightarrow$ The primal problem has a strong dual

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

$$v_P = v_D$$



## Discrete Empirical $\hat{P}_N$

$\Rightarrow$ $N$ samples of $z$

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

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

# Distributionally Robust Q-Learning

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



$s \rightarrow$  $Q(s,a_1)$
$Q(s,a_2)$
$\vdots$

## Bellman Operator $\mathcal{T}$

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

$\Rightarrow$ Contraction

$\underbrace{\sum_{a'\in A} \pi(a'|s) Q(s',a')}$

$\pi(a|s) = \dfrac{e^{Q(s,a)}}{\sum\limits_{a'\in A} e^{Q(s,a')}}$

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

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

$\Rightarrow$ worst case expected discounted returns

$\quad h(s') \Rightarrow$ 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 $\Rightarrow \delta = \underbrace{r_k + \gamma \max_a Q(s_{k+1}, a)}_{\text{targets}} - Q(s_k, a_k)$

$Q(s_k, a_k) \leftarrow Q(s_k, a_k) + \alpha \delta$

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

$\Rightarrow$ Double DQN
$\Rightarrow$ Duelling layer
$\Rightarrow$ Prioritized experience replay

# Solving DRO

## Lipschitz Approximation

* $\sup\limits_{P \in U} E_{s' \sim P}[-h(s')]$

$$\text{Lip}(\emptyset) = \sup\limits_{s \neq s'} \frac{|\emptyset(s) - \emptyset(s')|}{\|s - s'\|}$$

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

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

  LipSDP in python

  $\searrow$ neural network lipschitz const. estimation

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

$\Rightarrow$ Pessimistic upper bound

## Robust Program Approximation

$$V_P = V_D = \min\limits_{\lambda \geq 0} \left\{ \lambda \epsilon^p + \frac{1}{N} \sum_{i=1}^{N} \sup\limits_{z \in Z} \left[ -h(z) - \lambda \|z - z_i\|^p \right] \right\}$$

for $K > 0$

$$V_k := \sup\limits_{(z^{ik})_{ik} \in M_k} \frac{1}{NK} \sum_{i=1}^{N} \sum_{k=1}^{K} -h(z^{ik})$$

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

$$V_k \uparrow \sup\limits_{P \in U} E_{s' \sim P}[-h(s')] \quad \text{as} \quad K \to \infty$$

$\Rightarrow$ Local optimums can be found

$\Rightarrow$ Must solve for each transition in memory

$\Rightarrow$ Will be slow

$\Rightarrow$ Optimistic lower bound