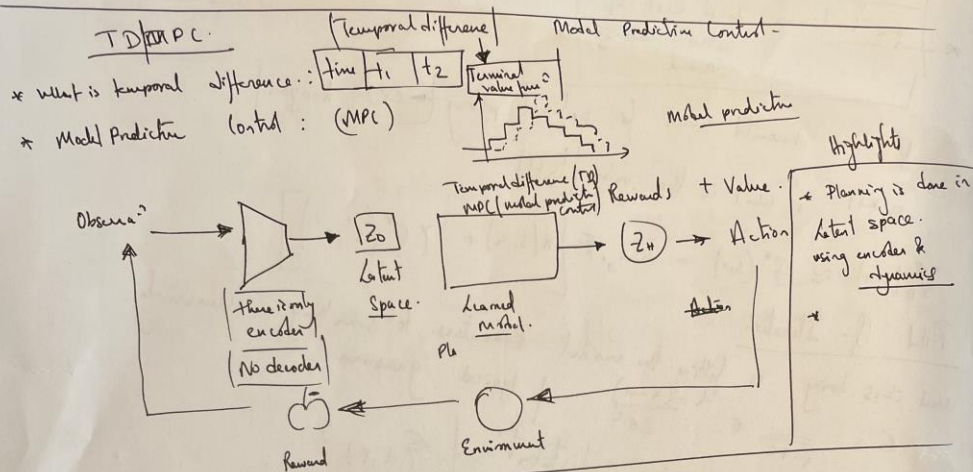


Notes:

Comparison of Planning & Learning algorithms
Compare these different models & architectures used.

APRIL 16/2025



MPC
Optimizing Based strategy \rightarrow time t

"complete optimization" \rightarrow Action 1 executed, $(x \times x)$ repeated. $(t+1)$

(model-free model-based environments)

Model Based Methods
Complex.
Worse Performance

Model-Free Methods
Simpler.
Better Performance

Fitted Q -iteration — long (infinite) Horizon step.

$$Q_0(s, a) \approx Q^*(s, a) = \max_{a'} \mathbb{E} \left[R(s, a) + \gamma Q^*(s', a') \right]$$

Expected return \downarrow Reward based on (s, a) (s', a') (state, action value step next)

$$\arg \min_{\theta} \mathbb{E}_{s, a \sim \pi} \| Q_{\theta}(s, a) - y \|_2$$

using values from the replay buffer
- just Expected Reward \mathbb{E} what is this.

where $Q^* = S \times A \Rightarrow \mathbb{R}$
 $\approx \mathbb{R}$ Reward

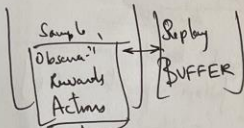
x contrastive models
vs generative models

APRIL 17/2025

Model of env.
↓ sample seq.
sample action

TD-learning
 $\phi_\theta(s, a)$

1. (learn model of environment) \rightarrow (2) sample (Action) | Evaluates trajectories using the learnt model.



enc

Latent space

\hat{a}_t
 \hat{r} ground truth reward / single step reward.
 q_t value TD learning.

↓

$\lambda(\phi) = \mathbb{E}_{(s, a, r, s')_{0:H}}$

$$\sim \theta \left[\sum_{t=0}^H \lambda^t \left(\| \hat{z}_t' - \text{sg}(\hat{h}(s_t')) \|_2^2 + \text{CE}(\hat{r}_t, r_t) + \text{CE}(\hat{q}_{t+1}, q_t) \right) \right]$$

predictor - ground truth reward

predict. one step to target q. target.

No decision

Online KL | Offline KL

high dimensional prob.

data \rightarrow policy - fast learning.

TD MPC-2

(100 different tasks) [100+ tasks]

Algorithm

Abbl. \rightarrow removing parts / specific
Observe \rightarrow how parts affect the system performance.

Plan
 \rightarrow Planning



Class 1:

TD-MPC

Key Concepts

- Recursive Dynamic Programming-Based

(CLASS-I)

Stochastic Reward

Example

finite Horizon: Limited no. of steps \rightarrow chooses actions accordingly.

Infinite Horizon: Unlimited no. of steps

at $\theta = 30^\circ$,

no. of sample = 20

Sample no. of Elites = 5 Best

recompute - new Gaussian distribution

20 random sample from distribution

narrow down the elite subset - close to the right answer.

- Find a good action trajectory - One finite Horizon.
- using Cross-Entropy Method. \rightarrow propose action trajectory
- sample action trajectory from multiple using Having way to propose what objective value would be?

Gaussian distribution

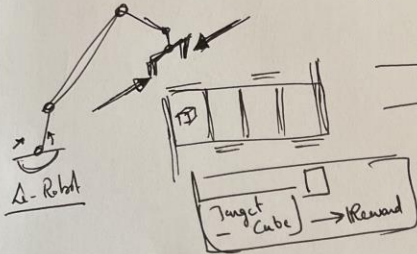
What objective value would be?

Predict IE R, find objective value from?

\rightarrow Plug & chug to Cross-Entropy Method

TD-MPC Class II

- Latent states \rightarrow faster \rightarrow lower dimensional state / lower processing speed.
- Focus on Relevant parts.
- Policy network predicts actions.
 - Gaussian takeoff
 - Cross-Entropy method.
- moving left is better, the robot can pick up cues on this.
- the imaginary bands would be useful.



- Practicing to push cube across
- TD learning before.

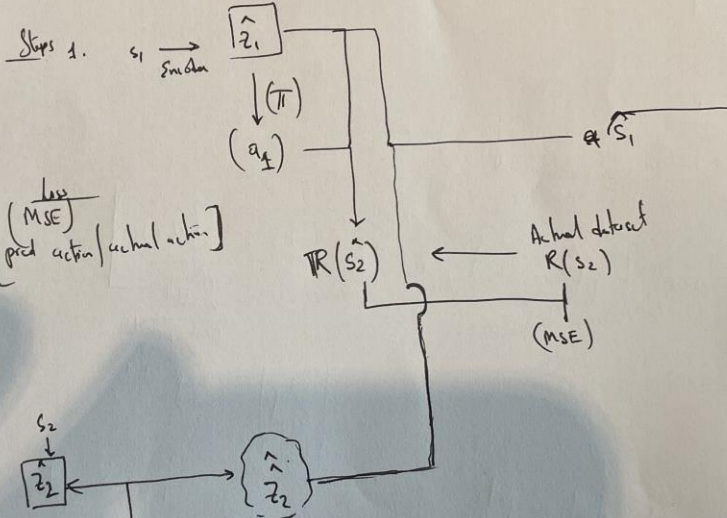
Online & offline learning
learn value func

\neq (Imitation learning) (more training needs to be done)

Parameter Horizon = 5

States	s_1	s_2	s_3	s_4	s_5	s_6
Actions	a_1	a_2	a_3	a_4	a_5	a_6
Rewards	R_1	R_2	R_3	R_4	R_5	R_6

Note.



MPC \rightarrow not optimal
 \rightarrow noise / Explor?

TD \rightarrow not optimal
due to noise

R^* Network $(s, a) \rightarrow IE$

$\rightarrow \max(Q_{\text{network}})$

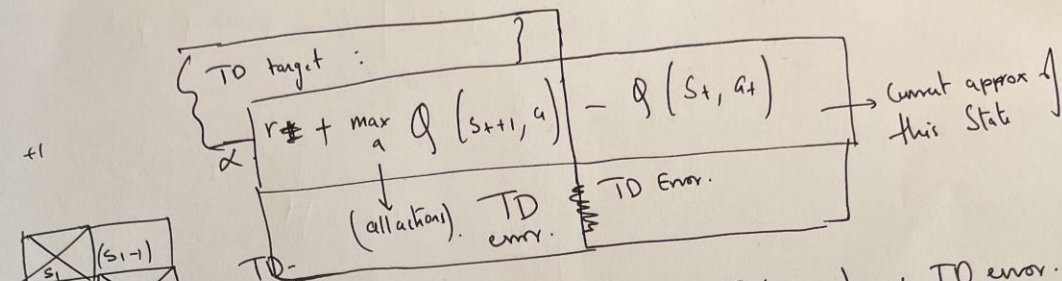
Advantage Weighted Action

(Q-Network)
(V-Network)

$\pi \rightarrow$ Action

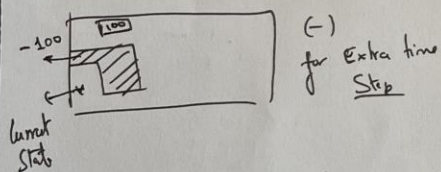
$R_{\text{net}} \rightarrow$ Reward Model

z_3
 z_4
 z_5
 z_6
 \rightarrow Consistency Loss



$$\text{Take } Q(s_t, a_t) \rightarrow \frac{Q(s_t, a_t)}{\text{Current Estimate}} \rightarrow \frac{Q(s_{t+1}, a_t)}{\text{Current Estimate}} + \text{TD error} + \alpha [\text{TD error}]$$

$$\begin{aligned} & (+100) \\ & -1 \\ & (+99) \\ & (s'_t + 100) \\ & (s_t - 1) \end{aligned}$$



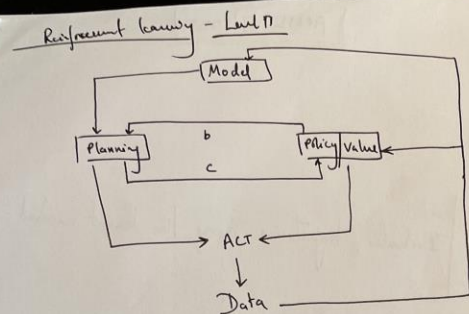
$$\text{MSE} \left[q(\hat{z}_1, a_1), r_2 + \frac{q(\hat{z}_2, \pi(\hat{z}_2))}{\text{policy model}} \right]$$

Best possible action in any state

Stop Grad here

not actual weights
actual exponential moving avg. \rightarrow stable target.

Reinforcement learning Overview:



Map Typical parameters:

MPP Typical parameters:
tuple: $\{S, A, T, R, p(s_0), \dots\}$

Envr. Transition func $T: S \times A \rightarrow \mathcal{P}(S)$
change from state to action

Reward free?
 $R: S \times A \times S \rightarrow IR$
 state 1 | Action | state 2
 derived state

Cannot find (s_t, a_t) .

Current state (s_t, a_t) .
 change. $s_t \sim T(\cdot | s_t, a_t)$ of transition func. $\left| \begin{array}{l} \text{However, biased, based upon condition} \\ \text{of initial/current state } s_t, \text{ \& follow up action } a_t \end{array} \right.$

Soln, $I_t = R(s_t, a_t, \underline{s_{t+1}})$

$$\text{trace of the env. } K \leq \sum_{k=0}^K r_{b+k}.$$

for $K = \infty$, infinite ~~time~~ horizon return.

discrete / continuous space.

Innov. is the algorithm.

$$g^{\pi}(s_t) = \mathbb{E}_{\pi, \gamma} \left[\sum_{k=0}^K \gamma^k r_{t+k} \mid s_t = s, a_t = a \right]$$

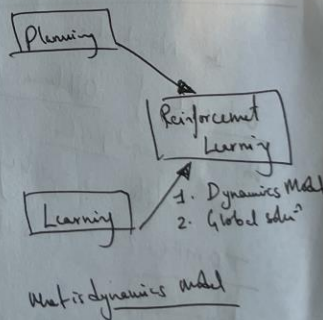
3:01 pm

c How is policy learnt?

- How is policy learnt?
- Or is the policy adopted / adapted from elsewhere?

Model Based RL

- learned Model & global soln.
- Known Model
- Planning over a learned model.



- Markov Decision problems (MDP)
probability of trajectory (state, action & Rewards) of length 'T':

APRIL 16 | 2025

$$P(Z) = P_0(s_0)$$

$$\text{Length } T: [T_1, T_2, T_3, \dots, T_T]$$

$$\text{Trajectory} = (T) \text{ for } T=1$$

States / example / Reward

States: s_0 start
 s_1 rewarded
 s_2 goal

Actions: a_{left}
 a_{right}

Rewards:
 $s_2 = r + 10$
others $r = 0$

$$P_0(s_0) = 1.0 \quad \text{Initial}$$

$$\pi(a_{\text{left}} | s_0) = 0.4 \quad \text{Policy}$$

Policy action (π): Deterministic
: Stochastic

Always $\frac{1.0}{0.8}$
0 ... 1.0
left right

π

Product Symbol (π) \rightarrow trajectory probability (chain sequentially)

$$P(Z) = P_0(s_0) \prod_{t=0}^{T-1} \pi(a_t | s_t) P_T(s_{t+1} | s_t, a_t) P_R(r_t | s_t, a_t, s_{t+1})$$

In a MDP

States (total $T+1$ states) $s_0, s_1, s_2, \dots, s_T$
Actions (total T actions) $a_0, a_1, a_2, \dots, a_{T-1}$
Rewards (total T rewards) r_0, r_1, \dots, r_{T-1}

$$P(Z) = P_0(s_0) \cdot \pi(a_0 | s_0) P_T(s_1 | s_0, a_0) P_R(r_0 | s_0, a_0, s_1) \cdot \pi(a_1 | s_1)$$

$$\times \pi(a_2 | s_2) P_T(s_3 | s_2, a_2) P_R(r_2 | s_2, a_2, s_3)$$

$$= 1.0 \cdot \left[\begin{array}{l} s_0, a_{\text{right}}, r_0=0, s_1, a_{\text{right}}, r_1=0, \\ s_2, a_{\text{right}}, r_2=10, s_3 \end{array} \right]$$

Probability transition (state $s_1 \rightarrow s_0$, action a_{right})

Initial

$$P_0(s_0) = 1.0$$

Policy

$$\pi(a_{\text{right}} | s_0) = 0.6$$

$$\pi(a_{\text{left}} | s_0) = 0.4$$

Transition

$$P_T(s_1 | s_0, a_{\text{right}}) = 0.8$$

$$P_T(s_2 | s_1, a_{\text{right}}) = 0.9$$

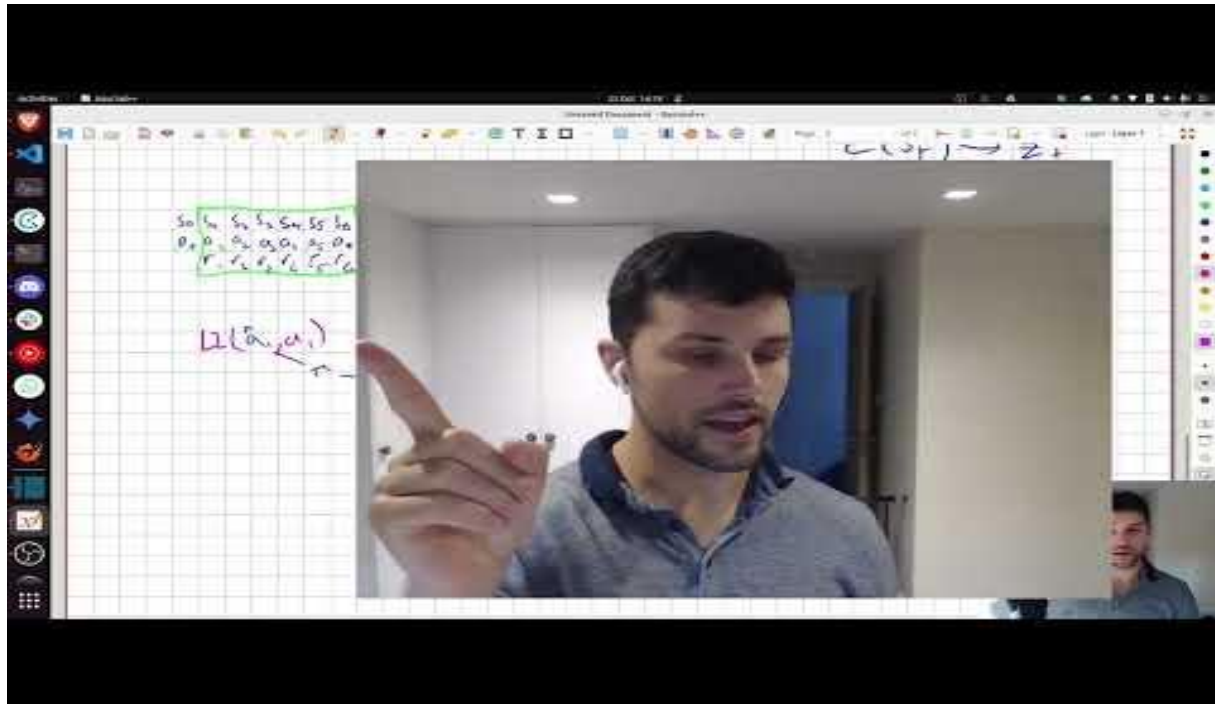
$$r=0, r=10$$

need to reach s_3

Sources:

Original Paper: <https://github.com/nicklashansen/tdmpc>

Reference learning: https://www.youtube.com/watch?v=_CKJJRAIvKI



<https://www.youtube.com/watch?v=--hDN4LLmPI>

Handwritten notes on a grid background, likely from a video lecture or presentation. The notes are written in black ink.

Top left: a_t with a horizontal arrow pointing right.

Top right: $reward(s_t, a_t) \rightarrow r_{t+1}$
 $dynamic(s_{t+1}, a_{t+1})$
 $q(s_{t+1}, a_{t+1}) \rightarrow g$

Center: A diagram of a rectangular environment. The top-left corner is shaded with diagonal lines. The top-left corner of the rectangle is labeled g . The bottom-left corner of the rectangle is labeled a .

Bottom right: A small video feed showing a person's face.