

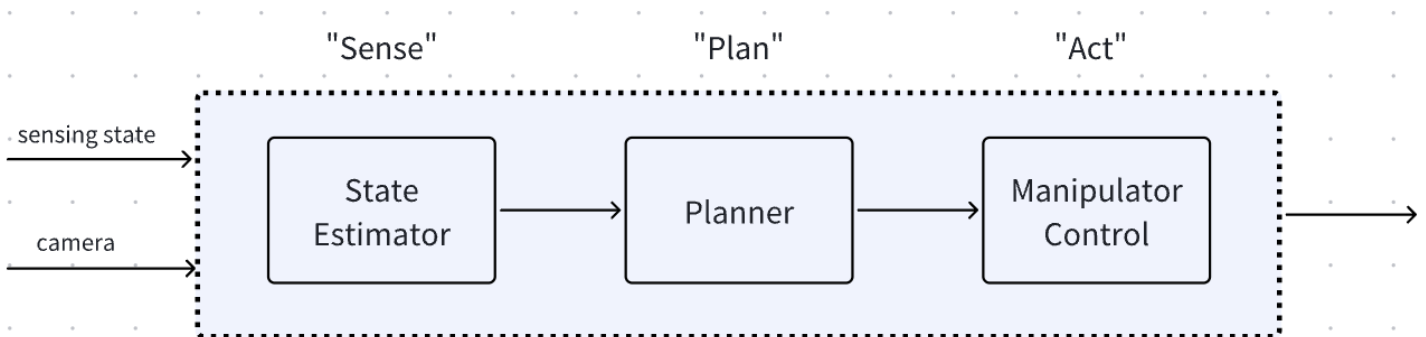
19. Reinforcement Learning 1

Key words

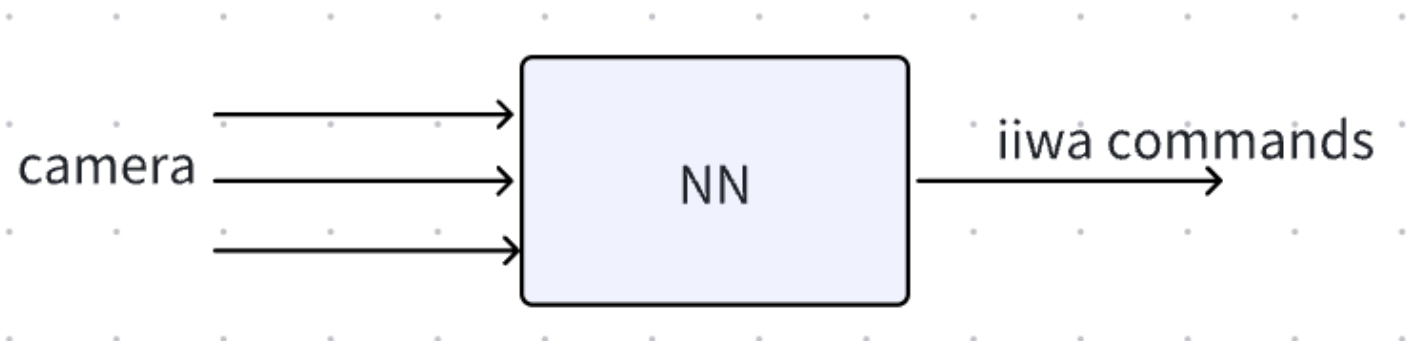
- Behavior cloning pipeline
- Foundation models
- Reinforcement Learning
 - model-free
 - model-based

Last Lecture, Visuomotor Policy

- How do we design controllers for visuomotor control?
- "model-based" control



- Behavior Cloning (Supervised sequence learning)



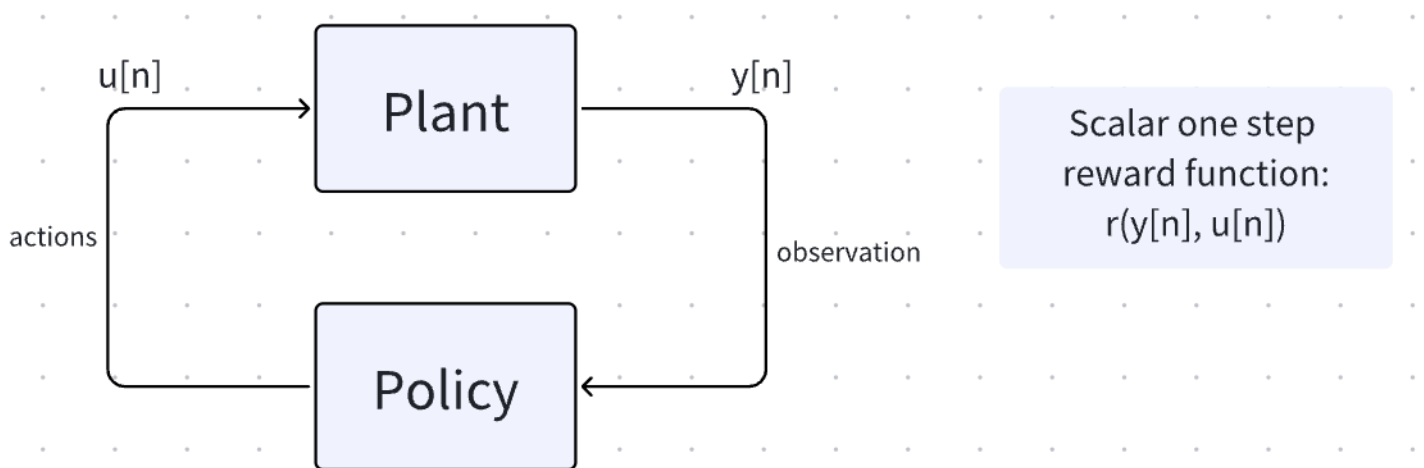
collect many demonstrations (e.g. human teleoperation)

- Reinforcement learning lies in the middle, requiring less supervision, but a harder optimization problem

- These methods are more similar than they are different: people asked, should i use BC, RL or model-based, as a matter of fact, if u get it done using one, u can also make it using the other methods. There are a lot connections.
- $y[n] = \pi_{\alpha}(y[n], u[n-1], y[n-1], \dots), \alpha : \text{NN parameters}$
-

Today's Topic outline

- Reinforcement Learning



- RL problem formulation
 - $\max_{\alpha} \sum_n r(y[n], u[n])$
 - subject to:
 - $x[n+1] = f(x[n], u[n]), x[0] = x_0$
 - $y[n] = g(x[n], u[n])$
 - $u[n] = \pi_{\alpha}(y[n], \dots u[n-1] \dots)$
- If I get a new robot, train from the beginning? Yes, most of the time, policy might be reused: [Cross embodiment](#)
- This is a general optimal control formulation, RL is a subset of Optimal Control, which emphasises
 - solution from trial & error
 - "black-box" optimization
 - stochastic updates, turn into probabilities, calculate the expectations

Handwritten mathematical equations on a chalkboard:

$$\max_{\alpha} E \left[\sum_{n=0}^{\infty} \gamma^n r(x[n], u[n]) \right]$$

$$P_x(x[n+1] | x[n], u[n])$$

$$P_y(y[n] | x[n], u[n])$$

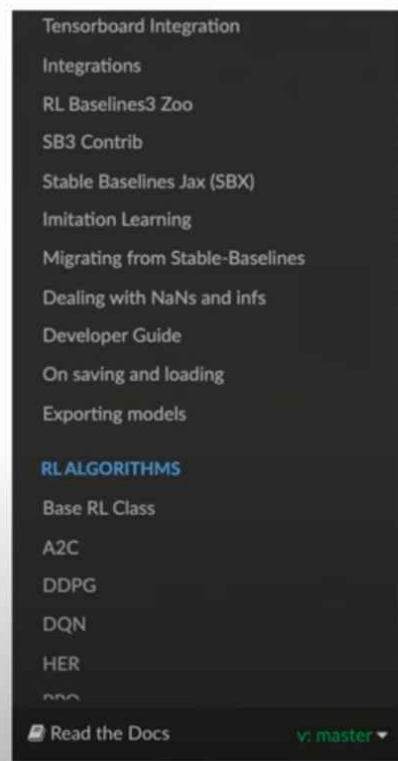
$$u[n] = \pi_{\alpha}(y[n], u[n-1], \dots)$$

$$P_{x_0}(x)$$

- neural net policies

1. RL is a harder optimization than BC

- cost depends on feedback loop through the world
- "delayed" reward
- collections of RL algorithms
 - RL in simulation use PPO(proximal Policy Optimization)
 - RL without simulation -> off-policy RL methods are popular (SAC)
 - Drake is not using GPU, not best choice for RL, Issac is. Drake is more focus on close the sim2real gap so its more accurate, Nvidia is not. So training from zero Use Issac, and then use Drake to train the second round.



multiprocessing.

Name	Box	Discrete	MultiDiscrete	MultiBinary	Multi Pro
ARS ¹	✓	✓	✗	✗	✓
A2C	✓	✓	✓	✓	✓
DDPG	✓	✗	✗	✗	✓
DQN	✗	✓	✗	✗	✓
HER	✓	✓	✗	✗	✓
PPO	✓	✓	✓	✓	✓
QR-DQN ¹	✗	✓	✗	✗	✓
RecurrentPPO ¹	✓	✓	✓	✓	✓
SAC	✓	✗	✗	✗	✓
TD3	✓	✗	✗	✗	✓
TQC ¹	✓	✗	✗	✗	✓
TRPO ¹	✓	✓	✓	✓	✓
Maskable PPO ¹	✗	✓	✓	✓	✓

[1] (1,2,3,4,5,6) Implemented in SB3 Contrib

2. Aside Parametrizing Dynamic Policies

- Linear dynamics example:

$x[n+1] = Ax[n] + Bu[n]$ $y[n] = Cx[n] + Dx[n]$	$y[n] = \sum_k \alpha_k y[n-k] + \sum_k \beta_k u[n-k]$
state space	ARX auto regressive model
recurrent neural network (LSTM)	GPT
Multibody Plant	Diffusion Policy

3. RL Recipe

1. Make the simulator
2. Write cost function, hardest part
3. Deep Policy gradient

4. Higher level messages

- Good software(gpu, pytorch etc)

- if you can write a simulation + elbow grease(cost function), most problems can be solved with RL
- Not RL vs BC vs Models, they all related

Ex: OpenAI Gym \Rightarrow Gymnasium

```

1 import gymnasium as gym
2
3 class FooEnv(gym.Env):
4     metadata = {'render.modes': ['human']}
5
6     def __init__(self):
7         ...
8     def step(self, action):
9         ...
10    def reset(self):
11        ...
12    def render(self, mode='human'):
13        ...
14    def close(self):
15        ...

```

```

1 import pydrake.all
2
3
4 builder = DiagramBuilder()
5 ....
6 diagram = builder.Build()
7 simulator = Simulator(diagram)
8
9
10 simulator.AdvanceTo(...)
11 observation = sensor_output_port->Eval(context)
12 reward = reward_output_port->Eval(context)
13
14 context = diagram.CreateDefaultContext()
15
16
17
18 meshcat.Publish(context)

```

5. Gym vs Drake

- Gym
 - no Context (-> no introspection, no deterministic playback)
 - floats only (no autodiff, no symbolic)
 - no access to multibody quantities (e.g., jacobians, inverse dynamics)
 - but can model anything
 - RL philosophy: "black box"; Drake: "glass box"



FINGER PIVOTING



SLIDING



FINGER GAITING

OpenAI - Learning Dexterity

"PPO has become the default reinforcement learning algorithm at OpenAI because of its ease of use and good performance."

<https://openai.com/blog/openai-baselines-ppo/>

Policy Architecture

Network

```
model = PPO('MlpPolicy', env, verbose=1, tensorboard_log=log)

stable_baselines3/common/policies.py#L435-L440

# Default network architecture, from stable-baselines
net_arch = [dict(pi=[64, 64], vf=[64, 64])]
```

approximately:

Actions

```
builder.ExportOutput(inv_dynamics.get_desired_position(), "action" COPY)
```

Observations

```
builder.ExportOutput(plant.get_state_output_port(), "observations")
```

Cost Function

```
1 angle_from_vertical = (box_state[2] % np.pi) - np.pi / 2 COPY
2 cost = 2 * angle_from_vertical**2 # box angle
3 cost += 0.1 * box_state[5]**2 # box velocity
4 effort = actions - finger_state[:2]
5 cost += 0.1 * effort.dot(effort) # effort
6 # finger velocity
7 cost += 0.1 * finger_state[2:].dot(finger_state[2:])
8 # Add 10 to make rewards positive (to avoid rewarding simulator
9 # crashes).
10 output[0] = 10 - cost
```

- in RL, if you have partial derivatives, do you use them? No
- Then how to minimize cost function without using gradients (python toolbox, never grad)
- Simplest form of RL, idea: weight perturbation
 - stochastic gradient descent
 - Eval $f()$ twice:
 - $f(\alpha)$, $f(\alpha + w)$, w : small random noise
 - $\Delta \vec{\alpha} = -\eta[f(\alpha + w) - f(\alpha)]w$
 - learning rate: $E[\Delta \alpha] \propto -\frac{\partial f}{\partial \alpha}$
 - $var(\Delta(\alpha))$ might be bad, converging slowly, one thing you can do:
 - Change the place where you do the sampling:
 - adding noise to α , scale with number of params in network:
 - $u[n] = \pi_{\alpha}(y[n], u[n-1]..) + w[n]$, add noise to policy output
 - variance grows with # of outputs * # of timesteps