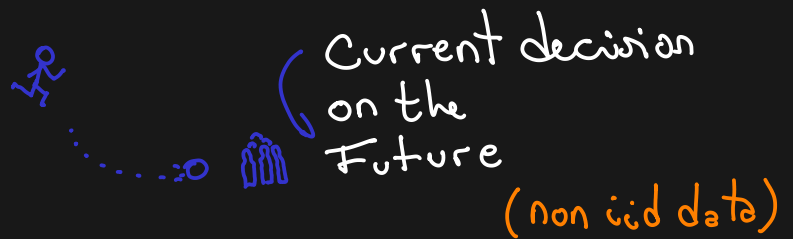# Reinforcement Learning

- Sequential decision making.
  - **Not** needed when:
    - ↳ isolated decision (classification, regression)
    - ↳ that decision does **not affect** future inputs or decisions.
  - Some applications can **not** ignore the dependence of the

Current decision
on the
Future

*(non iid data)*

The Plan:
1. Multi-task RL problem
2. Policy grads. & multi-task/meta counterparts/equivalents
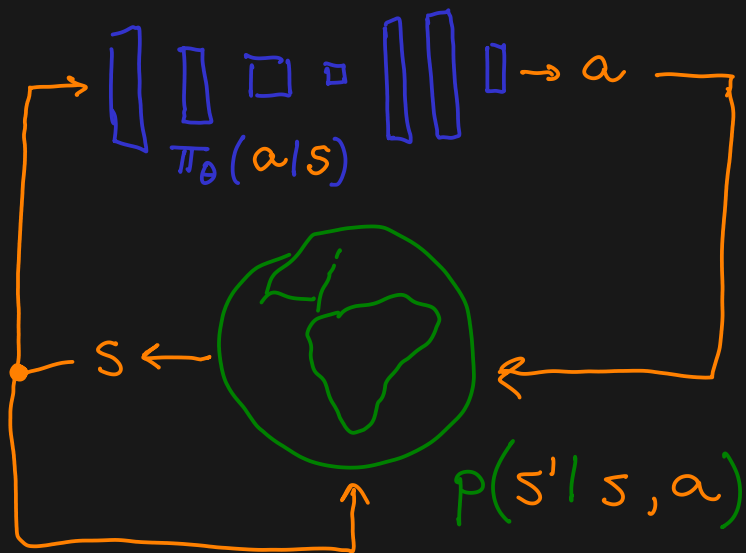3. Q-learning
4. Multi-task Q-learning.

1.

Imitation Learning
 ↳ Lots of data
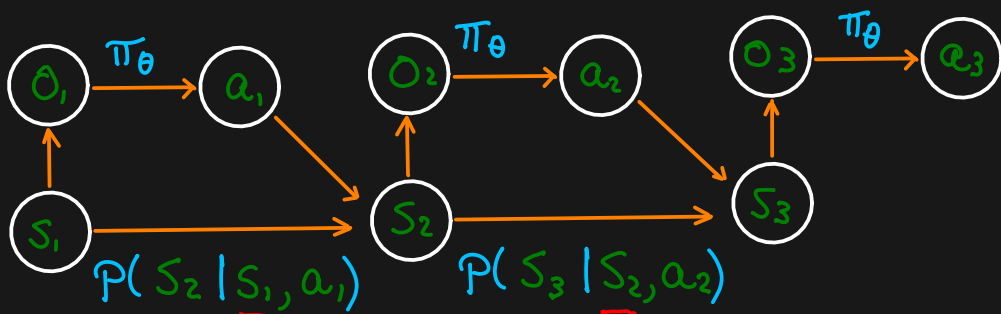 ↳ Short horizon decision problems
   (compound of errors otherwise)

Reward Function.

MDP: defined by
$$\begin{cases} \cdot \ s \\ \cdot \ a \\ \cdot \ r(s,a) \\ \cdot \ p(s'|s,a) \end{cases}$$

# The goal of RL:



## As Graphical Model:



$$P(S_{t+1} \mid S_t, a_t) : \text{Markov Prop.}$$

$$\theta^* = \arg\max \mathbb{E}_{(s,a) \sim p_\theta(s,a)} \left[ r(s,a) \right]$$

$$\theta^* = \arg\max \sum_{t=1}^{T} \mathbb{E}_{(s_t,a_t) \sim p_\theta(s_t,a_t)} \left[ r(s_t, a_t) \right]$$

Finite horizon: Episode = T steps

# RL Tasks

### Supervised Learn.

Task:

$$\mathcal{T}_i \triangleq \{ p_i(x), p_i(y|x), \mathcal{L}_i \}$$

$\underbrace{p_i(x)}_{\text{prior}} \quad \underbrace{p_i(y|x)}_{\text{likelihood}} \quad \underbrace{\mathcal{L}_i}_{\text{loss}}$

prior     likelihood     loss

$\underbrace{\phantom{p_i(x), p_i(y|x)}}_{\text{data generating distrib.}}$

### Reinforcement Learning

Task:

$$\mathcal{T}_i \triangleq \{ S_i, A_i, p_i(s_1), p_i(s'|s,a), r_i(s,a) \}$$

state   action   initial     dynamics       reward
space   space   state
           distribution

# Recommendation System as RL Tasks

each people = different task.

↳ personalized     $\left.\begin{array}{l} p_i(s'|s,a) \\ r_i(s,a) \end{array}\right\}$ Vary across tasks

# Character animation

- $r_i(s,a)$ vary across maneuvers



- $\left.\begin{array}{l} p_i(s_1) \\ p_i(s'|s,a) \end{array}\right\}$ vary across garment and initial state

- ## Multi-Robot RL

  $S_i, A_i, p_i(s_1), p_i(s'|s,a)$ vary across robots

Alternative view:
  A task identifier $z_i$ is part of the state $S$

$$s = (\bar{s}, z_i)$$
  ↑ original state

Going back to the first definition

$$\mathcal{T}_i \triangleq \left\{ S_i, A_i, p_i(s_1), p_i(s'|s,a), r_i(s,a) \right\}$$

We can see how as the task identifier is now part

of $s$ , $\begin{cases} p_i(s'|s,a), & \text{becomes} \\ r_i(s,a) \end{cases} \begin{cases} p(s'|s,a), \\ r(s,a) \end{cases}$

26/08/2020

$$\mathcal{T}_i \triangleq \left\{ S_i, A_i, p_i(s_1), p(s'|s,a), r(s,a) \right\}$$

- And now it can be cast as an MDP

$$\{\mathcal{T}_i\} = \left\{ \cup S_i, \cup A_i, \frac{1}{N}\sum_i p_i(s_1), p(s'|s,a), r(s,a) \right\}$$

This says: we can apply std. single task RL to the
multi-task problem with this view of multi-task RL

# The Goal of Multi-task

### Multi-task  ~~task~~ identifier

$$S = (\bar{s}, z_i)$$

$z_i$ could be: 1-hot task ID

lang. description

• Desired Goal State : $z_i = s_g$
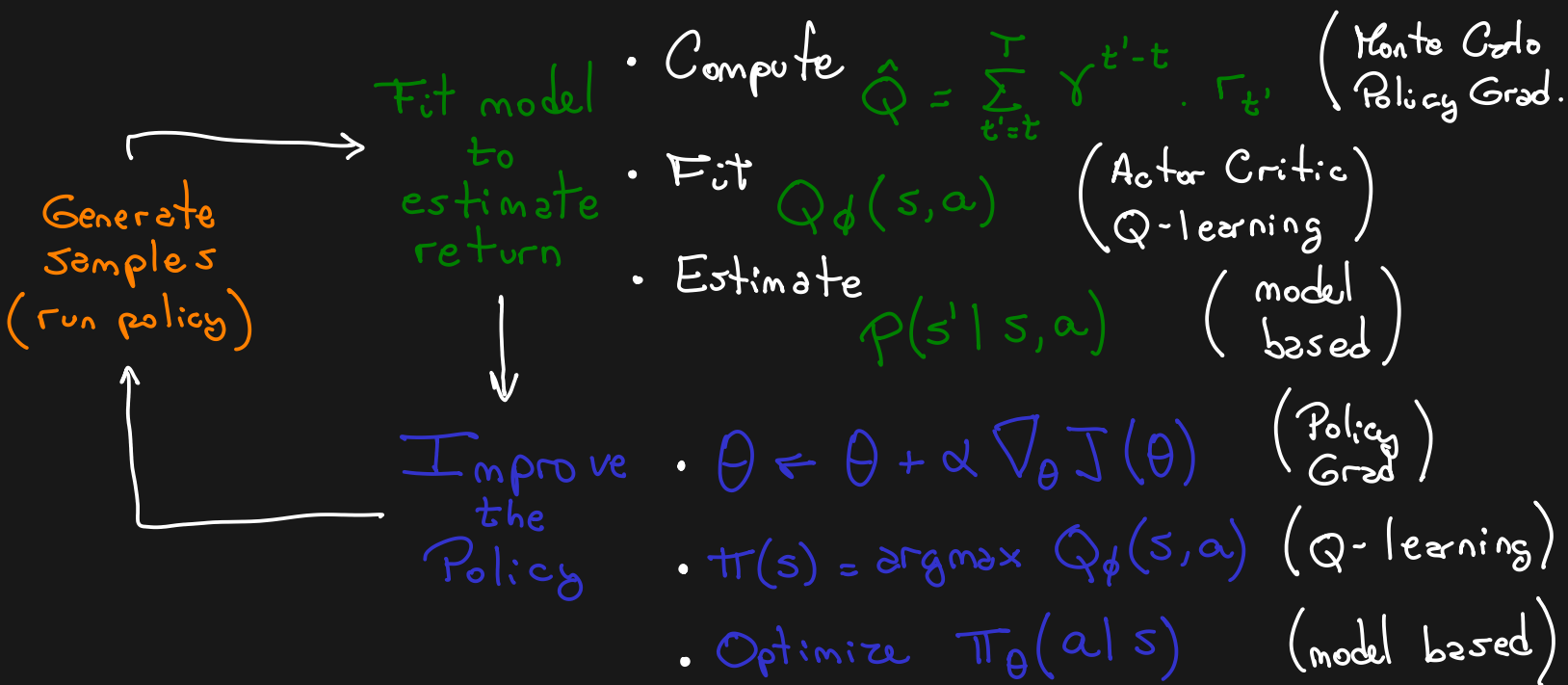
this is called

"Goal-Conditioned RL"

where the reward is some distance to $s_g$

$$r(s) = r(\bar{s}, s_g) = -d(\bar{s}, s_g)$$

eg: Euclidean $\ell_2$

Sparse 0/1 : $\mathbb{1}\{\bar{s} = s_g\}$

# RL Algorithm anatomy.

Generate samples (run policy)

Fit model to estimate return

• Compute $\hat{Q} = \sum_{t'=t}^{T} \gamma^{t'-t} \cdot r_{t'}$  (Monte Carlo Policy Grad.)

• Fit $Q_\phi(s,a)$  (Actor Critic) (Q-learning)

• Estimate $P(s'|s,a)$  (model based)

Improve the Policy

• $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$  (Policy Grad)

• $\pi(s) = \text{argmax}\ Q_\phi(s,a)$  (Q-learning)

• Optimize $\pi_\theta(a|s)$  (model based)

We focus to model-free for now.

# Evaluating the objective

$$\theta^* = \arg\max_\theta \; \mathbb{E}_{\tau \sim p_\theta(\tau)} \left[ \sum_t r(s_t, a_t) \right]$$

$$\underbrace{\qquad\qquad\qquad\qquad}_{J(\theta)}$$

$$J(\theta) \approx \frac{1}{N} \sum_i \sum_t r(s_{i,t}, a_{i,t})$$

sum over samples
from $\pi_\theta$

Samples from $\pi_\theta$

$r(s_{3,1}, a_{3,1})$

$t=0$   $t=1$

# Direct Policy Differentiation

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta(\tau)} \left[ r(\tau) \right] = \int \pi_\theta(\tau) \, r(\tau) \, d\tau$$

$\llcorner$ integrates over <u>trajectories</u>

dif:

$$\nabla_\theta J(\theta) = \int \underbrace{\nabla_\theta \pi_\theta(\tau)}_{= \pi_\theta(\tau) \cdot \frac{\nabla_\theta \pi_\theta(\tau)}{\pi_\theta(\tau)}} r(\tau) \, d\tau$$

$$= \int \pi_\theta(\tau) \cdot \nabla_\theta \log \pi_\theta(\tau) \cdot r(\tau) \cdot d\tau$$

$$= \mathbb{E}_{\tau \sim \pi_\theta(\tau)} \left[ \nabla_\theta \log \pi_\theta(\tau) \cdot r(\tau) \right]$$

# Probability of trajectory under policy $\pi$

$$\underbrace{\pi_\theta(s_1, a_1, \ldots, s_T, a_T)}_{\pi_\theta(\tau)} = p(s_1) \cdot \overset{T}{\underset{t=1}{\prod}} \pi_\theta(a_t | s_t) \cdot p(s_{t+1} | s_t, a_t)$$

over all trajectory time-steps

↑ initial state density    ↑ policy prob.    ↑ dynamics probability

$$\log \pi_\theta(\tau) = \log p(s_1) + \sum_{t=1}^{T} \log \pi_\theta(a_t | s_t) + \log p(s_{t+1} | s_t, a_t)$$

$$\nabla_\theta \log \pi_\theta(\tau) = \oslash + \sum_{t=1}^{T} \nabla_\theta \log \pi_\theta(a_t | s_t) + \oslash$$

does not depends on $\theta$          does not depends on $\theta$

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta(\tau)} \left[ \left( \sum_{t=1}^{T} \nabla_\theta \log \pi_\theta(a_t | s_t) \right) \left( \sum_{t=1}^{T} r(s_t, a_t) \right) \right]$$

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left( \sum_{t=1}^{T} \nabla_\theta \log \pi_\theta(a_{i,t} | s_{i,t}) \right) \left( \sum_{t=1}^{T} r(s_{i,t}, a_{i,t}) \right)$$

generate samples          estimate return

and then

$$\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$$

improve the policy

# REINFORCE

1. Sample $\{\tau^i\}$ from $\pi_\theta(a_t|s_t)$

2. $\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_i \left( \sum_t \nabla_\theta \log \pi_\theta(a_t^i, s_t^i) \right) \left( \sum_t r(s_t^i, a_t^i) \right)$

3. $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

# Imitation Learning

$\hookrightarrow$ Max. Likelihood of Expert actions

- Policy gradient

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_i \left( \sum_t \nabla_\theta \log \pi_\theta(a_t^i, s_t^i) \right) \left( \sum_t r(s_t^i, a_t^i) \right)$$

- Maximum Likelihood

$$\nabla_\theta J_{ML}(\theta) \approx \frac{1}{N} \sum_i \left( \sum_t \nabla_\theta \log \pi_\theta(a_t^i, s_t^i) \right)$$

Maximizes the probability of actions that have high reward
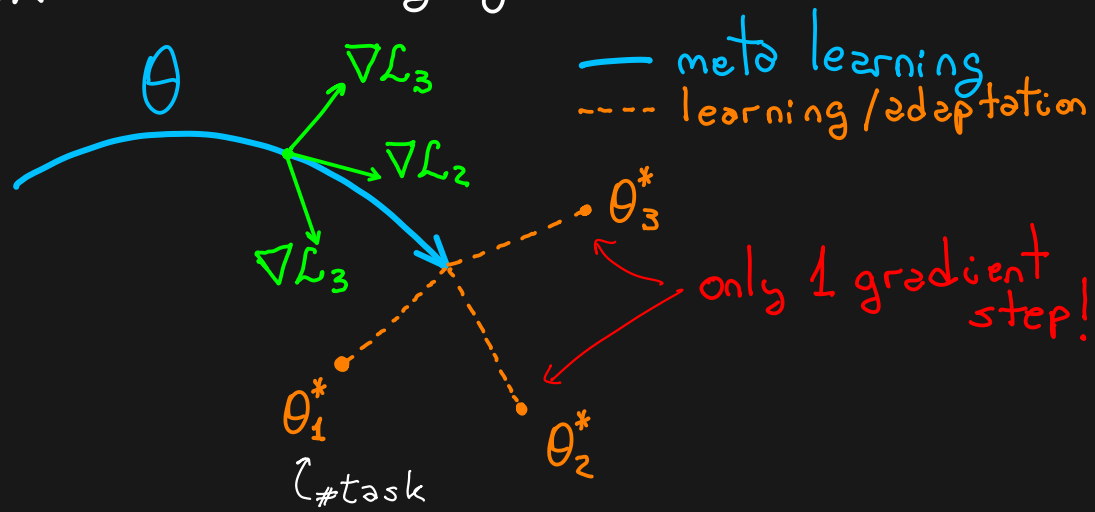
NOT weighted by reward

Policy gradient

Samples from $\pi_\theta$



best trajectory

$t=0$  $t=1$

We can directly apply RL + <u>Multi - task</u>

Example : MAML + Policy gradient



- $\theta$
- $\nabla L_3$
- $\nabla L_2$
- $\nabla L_3$
- $\theta_3^*$
- $\theta_1^*$
- $\theta_2^*$
- meta learning
- learning / adaptation
- only 1 gradient step!
- $\sum_{\#task}$

• <u>Black Box meta-learning</u> + Policy gradients

    3D Mazes - test on solve unseen mazes

• Combination with

    ⤷ Q-learning    } More difficult as are <u>not</u> gradient based
    ⤷ Actor Critic } algorithms : Bootstraping is Dynamic Programming

• Variance on Policy Grad.
    ⤷ can be mitigated with baselines, trust regions
        ⤷ see Hado van Hasselt video on baselines.

Importance weights can help to rehuse data.

# Value-Based RL

## Value function

- $V^{\pi}(s_t) = \sum_{t'=t}^{T} \mathbb{E}_{\pi} \left[ r(s_{t'}, a_{t'} | s_t) \right]$

## Q function

- $Q^{\pi}(s_t, a_t) = \sum_{t'=t}^{T} \mathbb{E}_{\pi} \left[ r(s_{t'}, a_{t'} | s_t, a_t) \right]$

equiv.

- $V^{\pi}(s_t) = \mathbb{E}_{a_t \sim \pi(\cdot | s_t)} \left[ Q^{\pi}(s_t, a_t) \right]$

## Bellman eq:

- For the optimal policy $\pi^*$

$$Q^*(s,a) = \mathbb{E}_{s' \sim p(\cdot | s, a)} \left[ r(s,a) + \gamma \max Q^*(s', a') \right]$$

↳ Best reward in the future

## Fitted Q-iteration: DP algorithm that leads to Bellman eq.

→ 1. collect dataset $\{(s_i, a_i, s_i', r_i)\}$

→ 2. set $y_i \leftarrow r(s_i, a_i) + \gamma \cdot \max_{a_i'} Q_\phi(s_i', a_i')$

K iters

3. set $\phi \leftarrow \arg\min_\phi \frac{1}{2} \sum_i \| Q_\phi(s_i, a_i) - y_i \|^2$

- Now we have $Q_\phi(s,a)$, get policy $\pi(a|s)$ from $\arg\max_a Q_\phi(s,a)$

$\int$ Off policy algo.
  ↳ can use replay buffer

- <u>Not</u> grad. descent algorithm : is a DP algo. ⟹ <span style="color:red">tricky to combine with MAML or Black Box approach.</span>
- Can be readily extended to multi-task / goal-conditioned RL

# Multi-task RL Algorithms

- Policy
  $\qquad$ ┌─ condition also on task $z_i$
  $$\pi_\theta(a|\bar{s}) \rightarrow \pi_\theta(a|\bar{s}, z_i)$$

- Q - function
  $$Q_\phi(\bar{s}, a) \rightarrow Q_\phi(\bar{s}, a, z_i)$$

Analogous as Multi task S.Learn., <u>but</u>
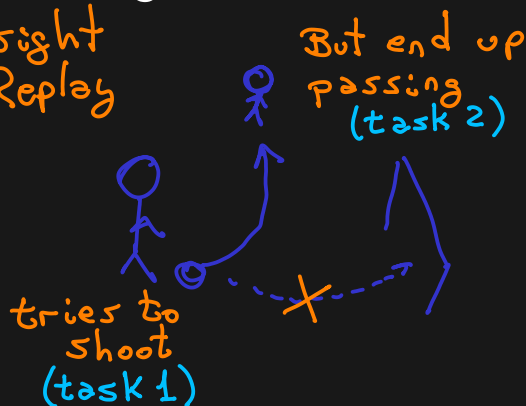
- Data distribution is controlled by the agent
- Data sharing across tasks?
- If known dynamics of MDP while changing across tasks.
  ↳ how can we leverage this knowledge?

(Spa: "Comprensión Retrospectiva")

# Hindsight relabeling / ake Hindsight Experience Replay "HER"

Relabel experience with other task id and store both

  ⟵ task 1 : { same experience }
  ↳ task 2 : { same experience }

aka Hindsight Experience Replay "HER"

But end up passing (task 2)

tries to shoot (task 1)

# Goal-conditioned RL with hindsight relabeling.

1. Collect data $\mathcal{D}_k = \{(s_{1:T}, a_{1:T}, s_{goal}, r_{1:T})\}$ using some policy.

2. Store data in replay buffer

$$\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_k$$

3. Perform hindsight relabeling

    a. Relabel experience in $\mathcal{D}_k$ using last state as goal.

$$\mathcal{D}_k' = \{(s_{1:T}, a_{1:T}, s_T, r_{1:T}')\}$$

    where $r_t' = -d(s_t, s_T)$ ← distance to (new) goal

    b. Store $\mathcal{D}_k'$ in Replay Buffer

$$\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_k'$$

4. Update Policy using replay buffer $\mathcal{D}$

Repeat $\uparrow$ 1.

Other relabeling in a.? "any state"

## Result:

    ↳ Exploration gets more efficient

    (exploring for 1 task helps exploration in other tasks)

# Multi-task RL with relabeling

Very similar structure:

1. Collect data $\mathcal{D}_k = \{(S_{1:T}, a_{1:T}, Z_i, r_{1:T})\}$ using some policy.

2. Store data in **replay buffer**

$$\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_k$$

3. Perform **hindsight relabeling**

> Choose
> - Randomly
> - task(s) in which the trajectory gets high reward.

   a. Relabel experience in $\mathcal{D}_k$ for task $\mathcal{T}_j$

$$\mathcal{D}'_k = \{(S_{1:T}, a_{1:T}, Z_j, r'_{1:T})\}$$

   where $r'_t = -r_j(S_t)$ ← return of state $S_t$ "following" task $\mathcal{T}_i$ }?

   b. Store $\mathcal{D}'_k$ in Replay Buffer

$$\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}'_k$$

## Requires

- Reward function form is known, evaluatable
- Dynamics consistent across goals / task
  - ↳ Relabeling is not that direct.

- Using off policy algorithm

Robot example
   much better with HER

# Image observations

   we need a distance function d between

$$r_t' = -d(S_t, S_T)$$

current state ← $S_t$
goal state ← $S_T$

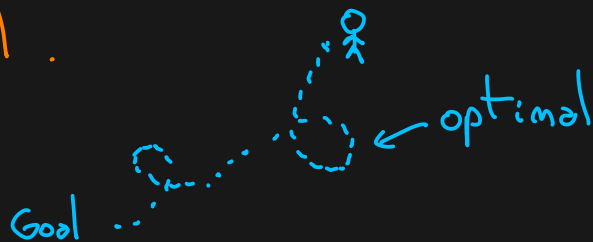Binary reward : $\begin{cases} 1 & \text{is same} \\ 0 & \text{if not} \end{cases}$

   ○ Sparse
   ○ Accurate

Random, unlabeled interaction
   ↳ optimal under the 0/1 reward of reaching last state.
   • We don't care what happens before the terminal state,
      all trajectories are optimal if they reach the
      desired image-goal.

"we have an optimal sample
for reaching the goal"

Goal ← optimal

If data is optimal
   ↳ use supervised imitation learning.

Similar structure again:

1. Collect data $\mathcal{D}_k = \{(S_{1:T}, a_{1:T}, S_T, r_{1:T})\}$ using some policy.

· ~~Store data in replay buffer~~

~~$\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_k$~~

2. Perform hindsight relabeling

    a. Relabel experience in $\mathcal{D}_k$ using last state $S_T$ as goal

$$\mathcal{D}_k' = \{(S_{1:T}, a_{1:T}, S_T, r'_{1:T})\} \leftarrow \text{relabeled data}$$
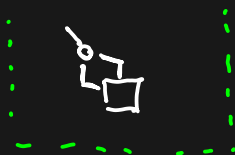
    where $r'_t = -d(S_t, S_T)$

    b. Store $\mathcal{D}_k'$ in Replay Buffer

$$\mathcal{D} \leftarrow \mathcal{D} \cup \mathcal{D}_k'$$

3. Update policy using <u>Supervised Imitation Learning</u>

on Replay Buffer $\mathcal{D}$

Example of robot w/ data from human: "Learn to reach image of a goal."

    ↳ Goal Image    vs    Current Image

↳ it works ⌐("/)⌐

Use insight to learn a better goal representation

1. Collect random, unlabeled iteraction data $\{(s_1, a_1, \ldots, a_{t-1}, s_t)\}$

2. Train a latent state representation
$$s \rightarrow x$$
and latent state model
$$f(x' | x, a)$$
st. if we plan a sequence of actions wrt goal state $s_t$
we recover the observed action sequence.

"This correspond to embedding a planer in latent space
into a goal condition policy"

3. Throw away latent space model

return goal representation X

("Distributional Planning Networks":

Accurate and shapped Reward Function
↳ not sparse as with $\mathbb{1}\{$ same image $\}$