# Imitation Learning

Why learn rewards?
→ to optimize what the person wants

otherwise: just copy the human
  $\checkmark$ → robot does things in human-like way
  $\checkmark$ → robot does things when we don't know $R/u$
  $\checkmark$ → robot has model of $\pi_H$, the human policy

$\xi_D \sim$ rollout from $\pi_D$    $s_D^0, a_D^0, s_D', a_D', \dots$

Behavio(al) Cloning: train $\pi$ s.t.  $\approx \pi_\theta(s) = \pi_D(s)$
  parametrize $\pi$ by $\theta$, e.g. ANN

$$\max_\theta \pi_\theta(a_D \mid s_D) \iff \max_\theta \sum_i \log \pi_\theta(a_D^i \mid s_D^i)$$

$$\min_\theta \mathbb{E}_{s \sim \pi_D} \left[ KL(\pi_D(a \mid s) \| \pi_\theta(a \mid s)) \right]$$

$$\iff \min_\theta \sum_i -\sum_a \pi_D(a \mid s_D^i) \cdot \log \left( \frac{\pi_\theta(a \mid s_D^i)}{\pi_D(a \mid s_D^i)} \right)$$

$\hookrightarrow \text{∫ in } a_D ?$    $\hookrightarrow \log \frac{x}{y} = \log x - \log y$

$$\approx \max_\theta \sum_i \log \pi_\theta(a_D^i \mid s_D^i)$$

what's wrong w. BC?

    ↳ assumes samples are IID

but we are in a sequential domain!

      ⇓

error accumulation

ALVIN (CMU):



new-ish state     even newer     off the road

theoretical argument:

supervised learning has chance of $\varepsilon$ error

if iid; T times: $T\varepsilon$ errors

BC is not iid: errors cumulate $\rightarrow T^2\varepsilon$

    ↳ once error $\rightarrow$ cost of $\perp$ @ each

                remaining timestp

$$\varepsilon \cdot (T-1) + \varepsilon(T-2) + \ldots + \varepsilon \cdot 1$$

fixes:

1) 2011 DAgger (dataset Aggregation)
→ • fit $\pi_\theta$ on $\mathcal{D}$
   • roll it out, collect induced states $S_{t=0}^T$
   • ask for action labels $a_{t=0}^T$
     <span style="color:magenta">(can people give you this?)</span>
   • $\mathcal{D} \leftarrow \mathcal{D} \cup (S, a)$

get "a policy" labels

(turns out injecting noise in the
demonstrator's actions to get
them to go off and recover is
probably enough (Dart '17)

2) practice with RL to "get back on" or
   "recover".
   2016 GAIL (generative adversarial inv. lear.)

think of IRL as
$$\max_C \min_\pi \left( \mathbb{E}_\pi[c(s,a)] - \lambda H(\pi) \right) - \mathbb{E}_{\pi_D}[c(s,a)]$$
w/out regularization of $c \to \pi$ matches demonst.
                              state-action occupancy

idea: search for a $\pi$ that does that more
                                        directly

$D(S,a) = 0$ if $(S,a)$ came from $\pi_D$, $1$ else

$\log D(S,a) = -\infty$ if $(S,a)$ came from $\pi_D$, $0$ else

train $\pi$ to minimize $D$; train $D$ to be low for $\pi_D$
but high otherwise:

$$\max_D \min_\pi \left( \mathbb{E}_\pi [\log D(S,a)] - \lambda H(\pi) \right) - \mathbb{E}_{\pi_D} [\log(1 - D(S,a))]$$

iterate between gradient on $D$
and R2 update on $\pi$