

Imitation and Inverse Reinforcement Learning

- Last time:
 - Turn-taking zero sum games
 - Markov Games
 - Incomplete Information Games
- Today:

Imitation and Inverse Reinforcement Learning

- **Last time:**
 - Turn-taking zero sum games
 - Markov Games
 - Incomplete Information Games
- **Today:**
 - What if you don't know the reward function and just want to act like an expert?
 - Imitation Learning
 - Inverse Reinforcement Learning

Trivia: When was the first car driven with a Neural Network?

Trivia: When was the first car driven with a Neural Network?



Dean Pomerleau
@deanpomerleau

Replies to @GTARobotics

GPU? Gez, ALVINN ran on 100 MFLOP CPU, ~10x slower than iWatch; Refrigerator-size & needed 5000 watt generator. [@olivercameron](#)

What's Hidden in the Hidden Layers?

The contents can be easy to find with a geometrical problem, but the hidden layers have yet to give up all their secrets

David S. Touretzky and Dean A. Pomerleau AUGUST 1989 • BY T E 231

tions, we fed the network road images taken under a wide variety of viewing angles and lighting conditions. It would be impractical to try to collect thousands of real road images for such a data set. Instead, we developed a synthetic road-image generator that can create as many training examples as we need.

To train the network, 1200 simulated road images are presented 40 times each, while the weights are adjusted using the back-propagation learning algorithm. This takes about 30 minutes on Carnegie Mellon's Warp systolic-array supercomputer. (This machine was designed at Carnegie Mellon and is built by General Electric. It has a peak rate of 100 million floating-point operations per second and can compute weight adjustments for back-propagation networks at a rate of 20 million connections per second.)

Once it is trained, ALVINN can accurately drive the NAVLAB vehicle at about 3½ miles per hour along a path through a wooded area adjoining the Carnegie Mellon campus, under a variety of weather and lighting conditions. This speed is nearly twice as fast as that achieved by non-neural-network algorithms running on the same vehicle. Part of the reason for this is that the forward pass of a back-propagation network can be computed quickly. It takes about 200 milliseconds on the Sun-3/160 workstation installed on the NAVLAB.

The hidden-layer representations ALVINN develops are interesting. When trained on roads of a fixed width, the net-

work chooses a representation in which hidden units act as detectors for complete roads at various positions and orientations. When trained on roads of variable

continued



Photo 1: The NAVLAB autonomous navigation test-bed vehicle and the road used for trial runs.

Trivia: When was the first car driven with a Neural Network?



Dean Pomerleau
@deanpomerleau

...

1995: 2797/2849 miles (98.2%)

Replies to @GTARobotics

GPU? Gez, ALVINN ran on 100 MFLOP CPU, ~10x slower than iWatch; Refrigerator-size & needed 5000 watt generator. [@olivercameron](#)

What's Hidden in the Hidden Layers?

The contents can be easy to find with a geometrical problem, but the hidden layers have yet to give up all their secrets

David S. Touretzky and Dean A. Pomerleau

AUGUST 1989 • BY T E 231

tions, we fed the network road images taken under a wide variety of viewing angles and lighting conditions. It would be impractical to try to collect thousands of real road images for such a data set. Instead, we developed a synthetic road-image generator that can create as many training examples as we need.

To train the network, 1200 simulated road images are presented 40 times each, while the weights are adjusted using the back-propagation learning algorithm. This takes about 30 minutes on Carnegie Mellon's Warp systolic-array supercomputer. (This machine was designed at Carnegie Mellon and is built by General Electric. It has a peak rate of 100 million floating-point operations per second and can compute weight adjustments for back-propagation networks at a rate of 20 million connections per second.)

Once it is trained, ALVINN can accurately drive the NAVLAB vehicle at about 3½ miles per hour along a path through a wooded area adjoining the Carnegie Mellon campus, under a variety of weather and lighting conditions. This speed is nearly twice as fast as that achieved by non-neural-network algorithms running on the same vehicle. Part of the reason for this is that the forward pass of a back-propagation network can be computed quickly. It takes about 200



Behavioral Cloning

Behavioral Cloning

$$\underset{\theta}{\text{maximize}} \prod_{(s,a) \in D} \pi_{\theta}(a \mid s)$$

Behavioral Cloning

$$\underset{\theta}{\text{maximize}} \prod_{(s,a) \in D} \pi_{\theta}(a \mid s)$$

Problem: Cascading Errors

How did ALVINN do it?

3.2. TRAINING "ON-THE-FLY" WITH REAL DATA

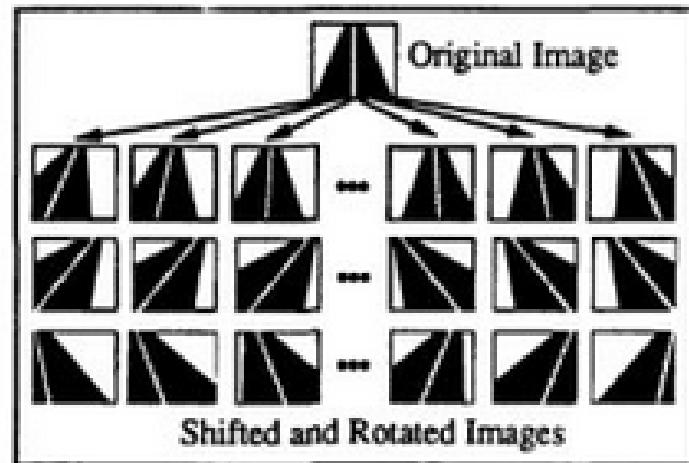


Figure 3.4: The single original video image is shifted and rotated to create multiple training exemplars in which the vehicle appears to be at different locations relative to the road.

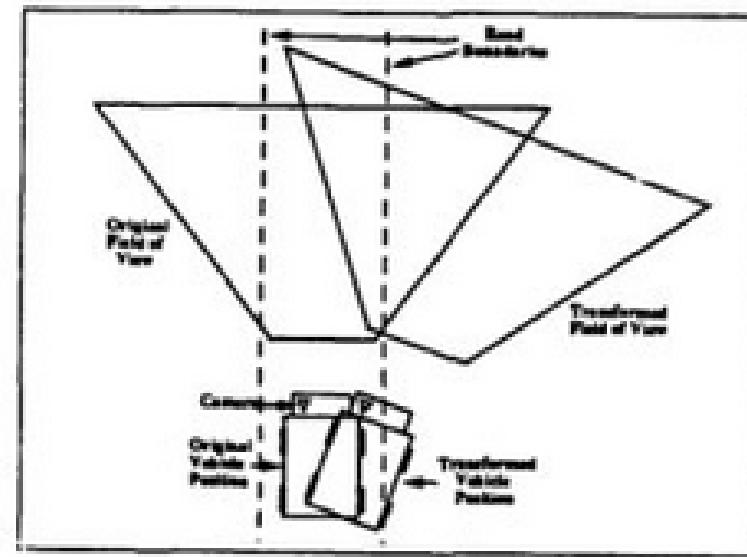


Figure 3.5: An aerial view of the vehicle at two different positions, with the corresponding sensor fields of view. To simulate the image transformation that would result from such a change in position and orientation of the vehicle, the overlap between the two field of view trapezoids is computed and used to direct resampling of the original image.

How did NVIDIA do it in 2016?

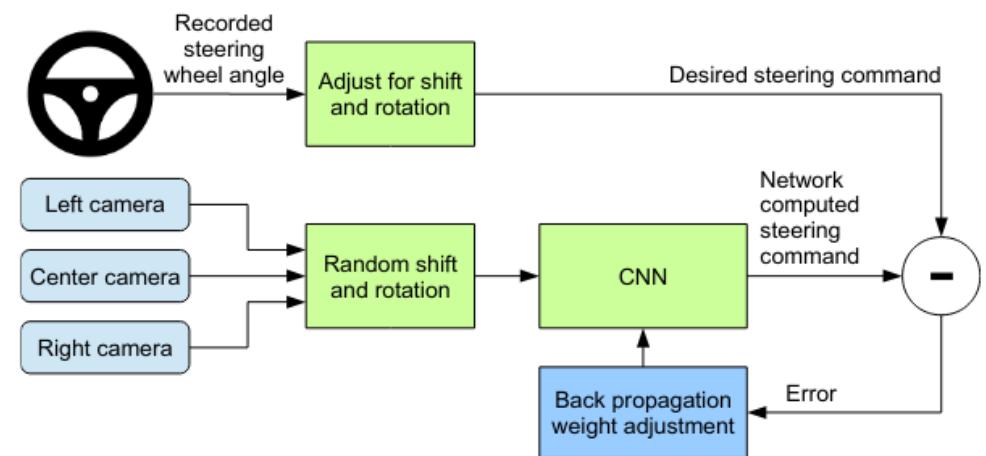
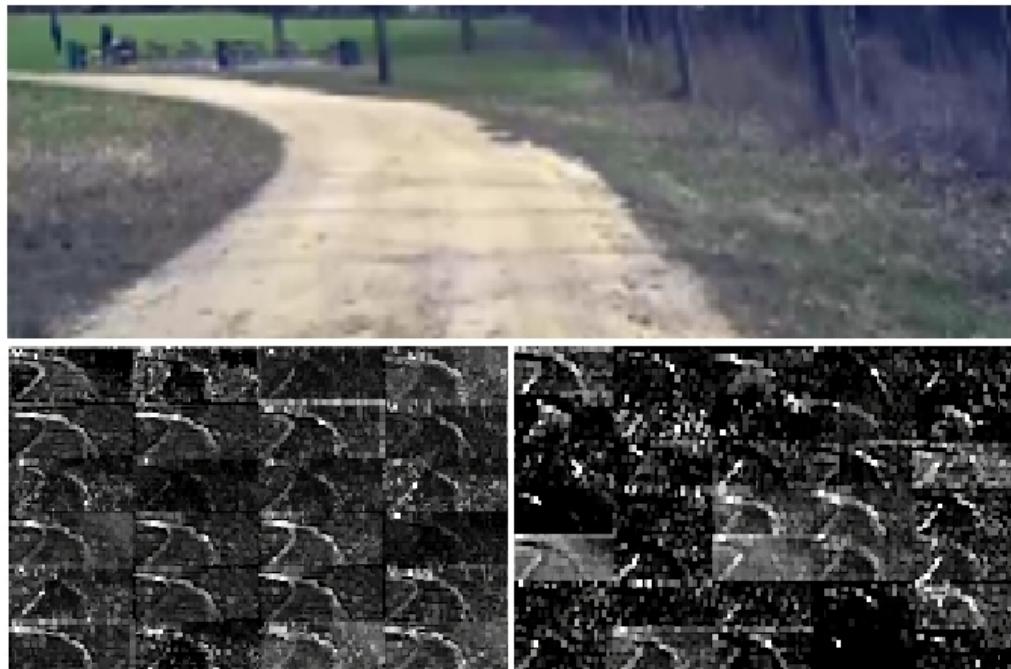


Figure 2: Training the neural network.

Dataset Aggregation (DAgger)

```
function optimize(M::DataSetAggregation, D, θ)
    ℙ, bc, k_max, m = M.ℙ, M.bc, M.k_max, M.m
    d, b, πE, πθ = M.d, M.b, M.πE, M.πθ
    θ = optimize(bc, D, θ)
    for k in 2:k_max
        for i in 1:m
            s = rand(b)
            for j in 1:d
                push!(D, (s, πE(s)))
                a = rand(πθ(θ, s))
                s = rand(ℙ.T(s, a))
            end
        end
        θ = optimize(bc, D, θ)
    end
    return θ
end
```

Dataset Aggregation (DAgger)

```
function optimize(M::DataSetAggregation, D, θ)
    ℙ, bc, k_max, m = M.ℙ, M.bc, M.k_max, M.m
    d, b, πE, πθ = M.d, M.b, M.πE, M.πθ
    θ = optimize(bc, D, θ)
    for k in 2:k_max
        for i in 1:m
            s = rand(b)
            for j in 1:d
                push!(D, (s, πE(s)))
                a = rand(πθ(θ, s))
                s = rand(ℙ.T(s, a))
            end
        end
        θ = optimize(bc, D, θ)
    end
    return θ
end
```

} rollout

Dataset Aggregation (DAgger)

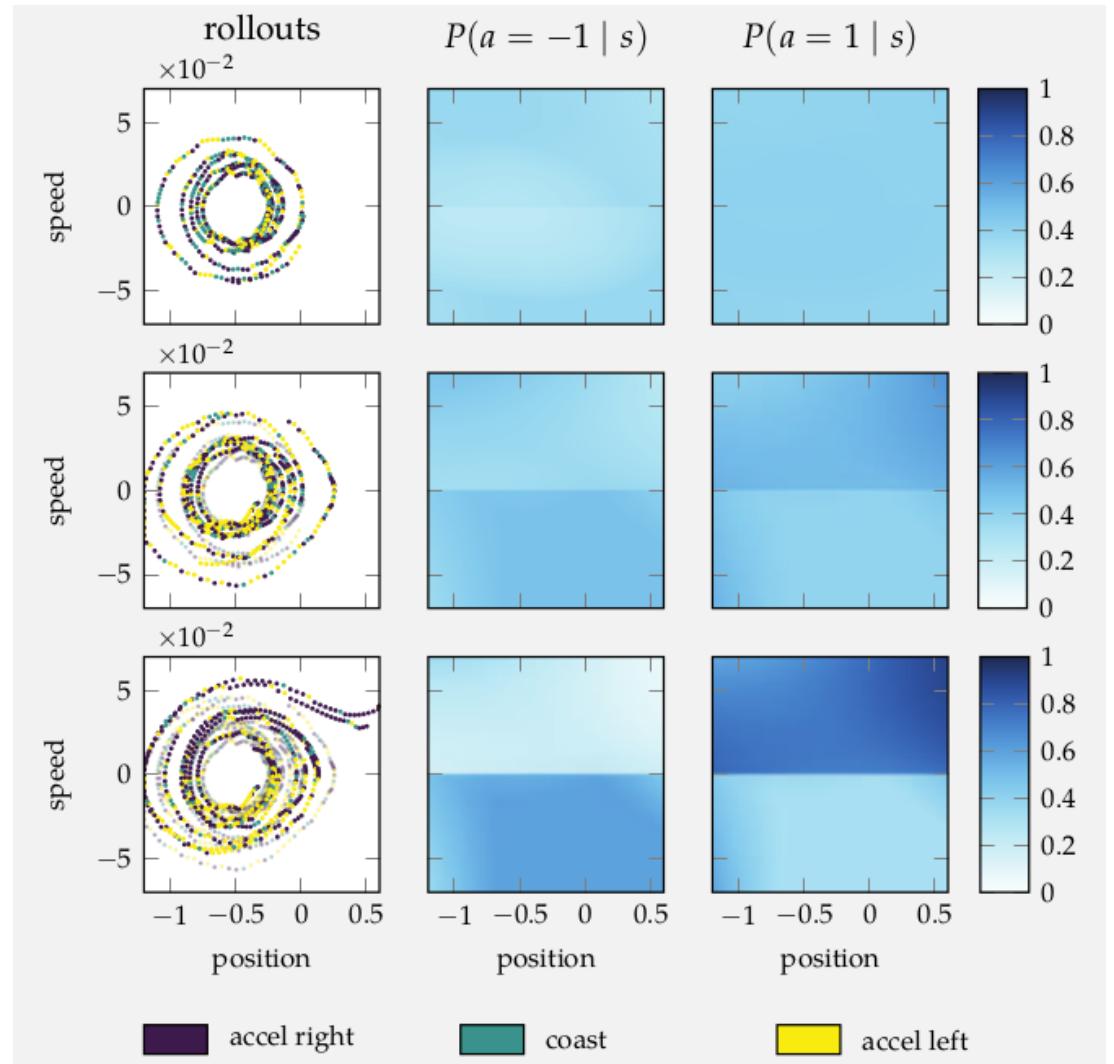
```
function optimize(M::DataSetAggregation, D, θ)
    ℙ, bc, k_max, m = M.ℙ, M.bc, M.k_max, M.m
    d, b, πE, πθ = M.d, M.b, M.πE, M.πθ
    θ = optimize(bc, D, θ)
    for k in 2:k_max
        for i in 1:m
            s = rand(b)
            for j in 1:d
                push!(D, (s, πE(s)))
                a = rand(πθ(θ, s))
                s = rand(ℙ.T(s, a))
            end
        end
        θ = optimize(bc, D, θ)
    end
    return θ
end
```

Gather from expert
} rollout

Dataset Aggregation (DAgger)

```
function optimize(M::DataSetAggregation, D, θ)
    P, bc, k_max, m = M.P, M.bc, M.k_max, M.m
    d, b, πE, πθ = M.d, M.b, M.πE, M.πθ
    θ = optimize(bc, D, θ)
    for k in 2:k_max
        for i in 1:m
            s = rand(b)
            for j in 1:d
                push!(D, (s, πE(s)))
                a = rand(πθ(θ, s))
                s = rand(P.T(s, a))
            end
        end
        θ = optimize(bc, D, θ)
    end
    return θ
end
```

Gather from expert
} rollout



Stochastic Mixing Iterative Learning (SMILe)

```
function optimize(M::SMILe, θ)
    P, bc, k_max, m = M.P, M.bc, M.k_max, M.m
    d, b, β, πE, πθ = M.d, M.b, M.β, M.πE, M.πθ
    A, T = P.A, P.T
    θs = []
    π = s → πE(s)
    for k in 1:k_max
        # execute latest π to get new data set D
        D = []
        for i in 1:m
            s = rand(b)
            for j in 1:d
                push!(D, (s, πE(s)))
                a = π(s)
                s = rand(T(s, a))
            end
        end
        # train new policy classifier
        θ = optimize(bc, D, θ)
        push!(θs, θ)
        # compute a new policy mixture
        Pπ = Categorical(normalize([(1-β)^(i-1) for i in 1:k], 1))
        π = s → begin
            if rand() < (1-β)^(k-1)
                return πE(s)
            else
                return rand(Categorical(πθ(θs[rand(Pπ)], s)))
            end
        end
    end
    Ps = normalize([(1-β)^(i-1) for i in 1:k_max], 1)
    return Ps, θs
end
```

Stochastic Mixing Iterative Learning (SMILe)

```
function optimize(M::SMILe, θ)
    P, bc, k_max, m = M.P, M.bc, M.k_max, M.m
    d, b, β, πE, πθ = M.d, M.b, M.β, M.πE, M.πθ
    A, T = P.A, P.T
    θs = []
    π = s → πE(s)
    for k in 1:k_max
        # execute latest π to get new data set D
        D = [] ← reset D
        for i in 1:m
            s = rand(b)
            for j in 1:d
                push!(D, (s, πE(s)))
                a = π(s)
                s = rand(T(s, a))
            end
        end
        # train new policy classifier
        θ = optimize(bc, D, θ)
        push!(θs, θ)
        # compute a new policy mixture
        Pπ = Categorical(normalize([(1-β)^(i-1) for i in 1:k], 1))
        π = s → begin
            if rand() < (1-β)^(k-1)
                return πE(s)
            else
                return rand(Categorical(πθ(θs[rand(Pπ)], s)))
            end
        end
    end
    Ps = normalize([(1-β)^(i-1) for i in 1:k_max], 1)
    return Ps, θs
end
```

Stochastic Mixing Iterative Learning (SMILe)

```
function optimize(M::SMILe, θ)
    P, bc, k_max, m = M.P, M.bc, M.k_max, M.m
    d, b, β, πE, πθ = M.d, M.b, M.β, M.πE, M.πθ
    A, T = P.A, P.T
    θs = []
    π = s → πE(s)
    for k in 1:k_max
        # execute latest π to get new data set D
        D = [] ← reset D
        for i in 1:m
            s = rand(b)
            for j in 1:d
                push!(D, (s, πE(s)))
                a = π(s)
                s = rand(T(s, a))
            end
        end
        # train new policy classifier
        θ = optimize(bc, D, θ)
        push!(θs, θ)
        # compute a new policy mixture
        Pπ = Categorical(normalize([(1-β)^(i-1) for i in 1:k], 1))
        π = s → begin
            if rand() < (1-β)^(k-1)
                return πE(s)
            else
                return rand(Categorical(πθ(θs[rand(Pπ)], s)))
            end
        end
    end
    Ps = normalize([(1-β)^(i-1) for i in 1:k_max], 1)
    return Ps, θs
end
```

} Gather
data

Stochastic Mixing Iterative Learning (SMILe)

```
function optimize(M::SMILe, θ)
    P, bc, k_max, m = M.P, M.bc, M.k_max, M.m
    d, b, β, πE, πθ = M.d, M.b, M.β, M.πE, M.πθ
    A, T = P.A, P.T
    θs = []
    π = s → πE(s)
    for k in 1:k_max
        # execute latest π to get new data set D
        D = [] ← reset D
        for i in 1:m
            s = rand(b)
            for j in 1:d
                push!(D, (s, πE(s))) } Gather
                a = π(s) data
                s = rand(T(s, a))
            end
        end
        # train new policy classifier
        θ = optimize(bc, D, θ) ← train only on D
        push!(θs, θ)
        # compute a new policy mixture
        Pπ = Categorical(normalize([(1-β)^(i-1) for i in 1:k], 1))
        π = s → begin
            if rand() < (1-β)^(k-1)
                return πE(s)
            else
                return rand(Categorical(πθ(θs[rand(Pπ)], s)))
            end
        end
    end
    Ps = normalize([(1-β)^(i-1) for i in 1:k_max], 1)
    return Ps, θs
end
```

Stochastic Mixing Iterative Learning (SMILe)

```
function optimize(M::SMILe, θ)
    P, bc, k_max, m = M.P, M.bc, M.k_max, M.m
    d, b, β, πE, πθ = M.d, M.b, M.β, M.πE, M.πθ
    A, T = P.A, P.T
    θs = []
    π = s → πE(s)
    for k in 1:k_max
        # execute latest π to get new data set D
        D = [] ← reset D
        for i in 1:m
            s = rand(b)
            for j in 1:d
                push!(D, (s, πE(s))) } Gather data
                a = π(s)
                s = rand(T(s, a))
            end
        end
        # train new policy classifier
        θ = optimize(bc, D, θ) ← train only on D
        push!(θs, θ)
        # compute a new policy mixture
        Pπ = Categorical(normalize([(1-β)^(i-1) for i in 1:k], 1))
        π = s → begin
            if rand() < (1-β)^(k-1)
                return πE(s)
            else
                return rand(Categorical(πθ(θs[rand(Pπ)], s)))
            end
        end
    end
    Ps = normalize([(1-β)^(i-1) for i in 1:k_max], 1)
    return Ps, θs
end
```

Stochastic Mixing Iterative Learning (SMILe)

```
function optimize(M::SMILe, θ)
    P, bc, k_max, m = M.P, M.bc, M.k_max, M.m
    d, b, β, πE, πθ = M.d, M.b, M.β, M.πE, M.πθ
    A, T = P.A, P.T
    θs = []
    π = s → πE(s)
    for k in 1:k_max
        # execute latest π to get new data set D
        D = [] ← reset D
        for i in 1:m
            s = rand(b)
            for j in 1:d
                push!(D, (s, πE(s))) } Gather data
                a = π(s)
                s = rand(T(s, a))
            end
        end
        # train new policy classifier
        θ = optimize(bc, D, θ) ← train only on D
        push!(θs, θ)
        # compute a new policy mixture
        Pπ = Categorical(normalize([(1-β)^(i-1) for i in 1:k], 1))
        π = s → begin
            if rand() < (1-β)^(k-1)
                return πE(s)
            else
                return rand(Categorical(πθ(θs[rand(Pπ)], s)))
            end
        end
        Ps = normalize([(1-β)^(i-1) for i in 1:k_max], 1)
    return Ps, θs
end
```

Mix Policies

$$(1 - \beta)^k$$

Stochastic Mixing Iterative Learning (SMILe)

```

function optimize(M::SMILe, θ)
    P, bc, k_max, m = M.P, M.bc, M.k_max, M.m
    d, b, β, πE, πθ = M.d, M.b, M.β, M.πE, M.πθ
    A, T = P.A, P.T
    θs = []
    π = s → πE(s)
    for k in 1:k_max
        # execute latest π to get new data set D
        D = [] ← reset D
        for i in 1:m
            s = rand(b)
            for j in 1:d
                push!(D, (s, πE(s)))
                a = π(s)
                s = rand(T(s, a))
            end
        end
        # train new policy classifier
        θ = optimize(bc, D, θ) ← train only on D
        push!(θs, θ)
        # compute a new policy mixture
        Pπ = Categorical(normalize([(1-β)^(i-1) for i in 1:k], 1))
        π = s → begin
            if rand() < (1-β)^(k-1)
                return πE(s)
            else
                return rand(Categorical(πθ(θs[rand(Pπ)], s)))
            end
        end
    end
    Ps = normalize([(1-β)^(i-1) for i in 1:k_max], 1)
    return Ps, θs
end

```

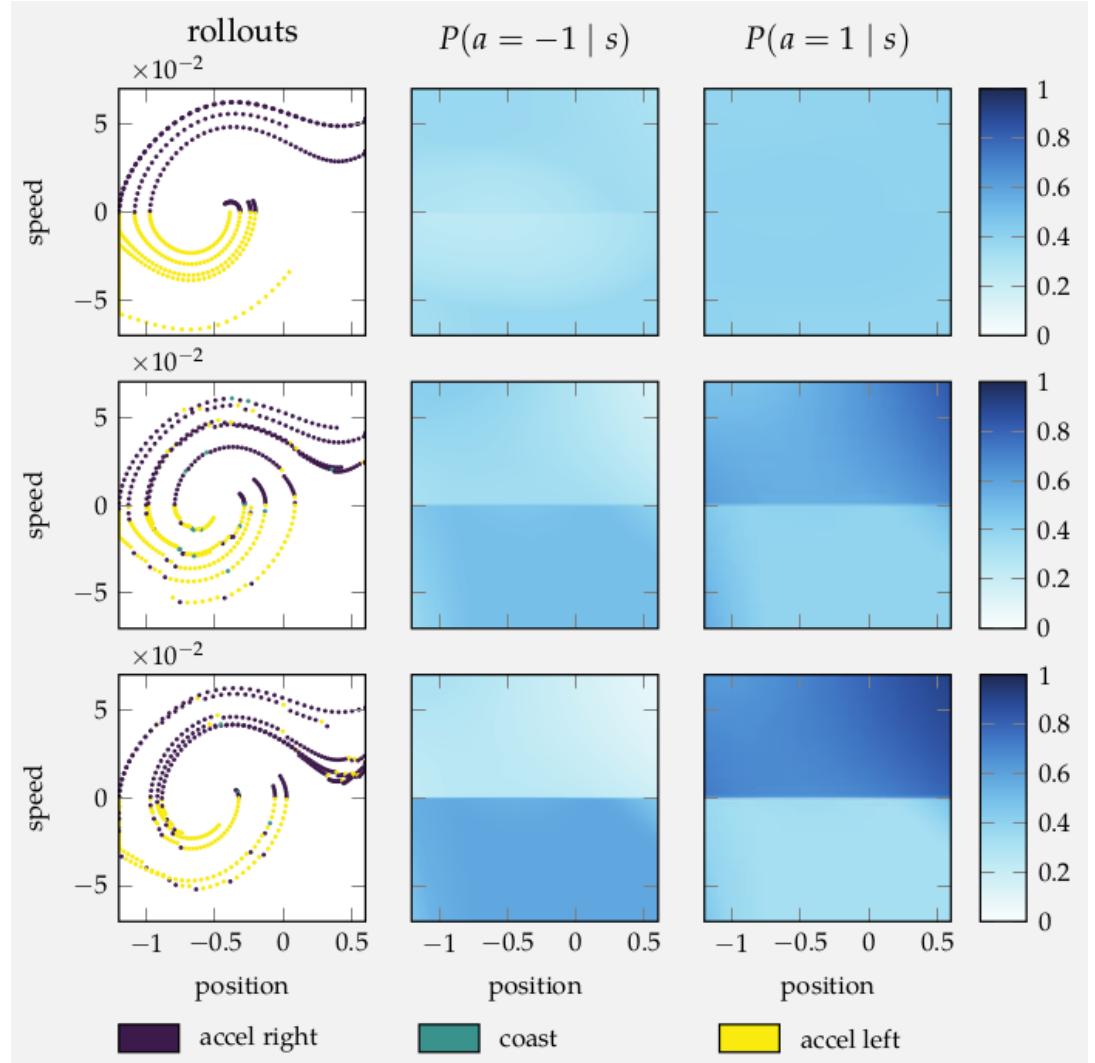
reset D

Gather data

train only on D

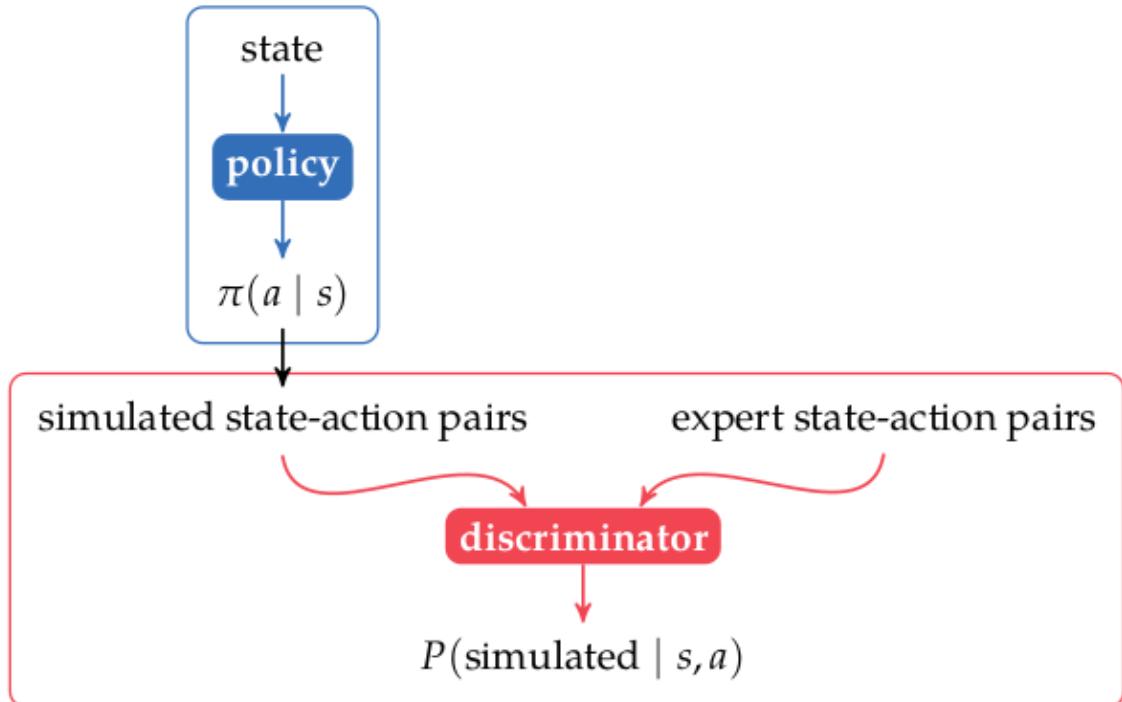
Mix Policies

$(1 - \beta)^k$

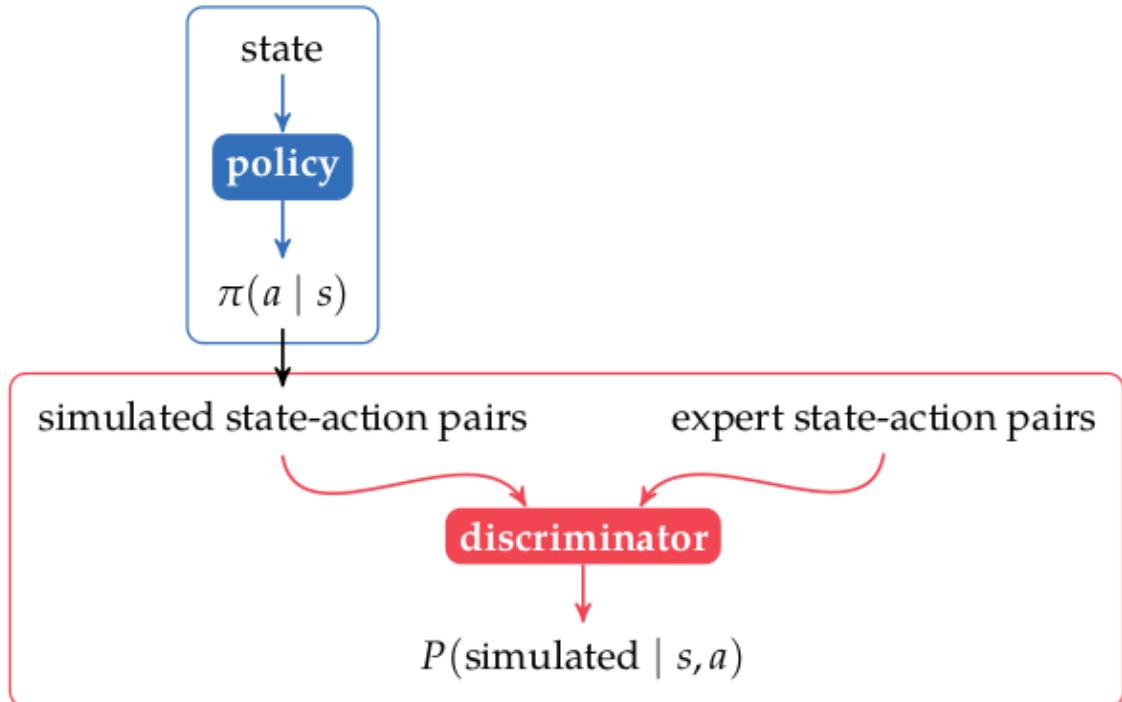


Generative Adversarial Imitation Learning (GAIL)

Generative Adversarial Imitation Learning (GAIL)

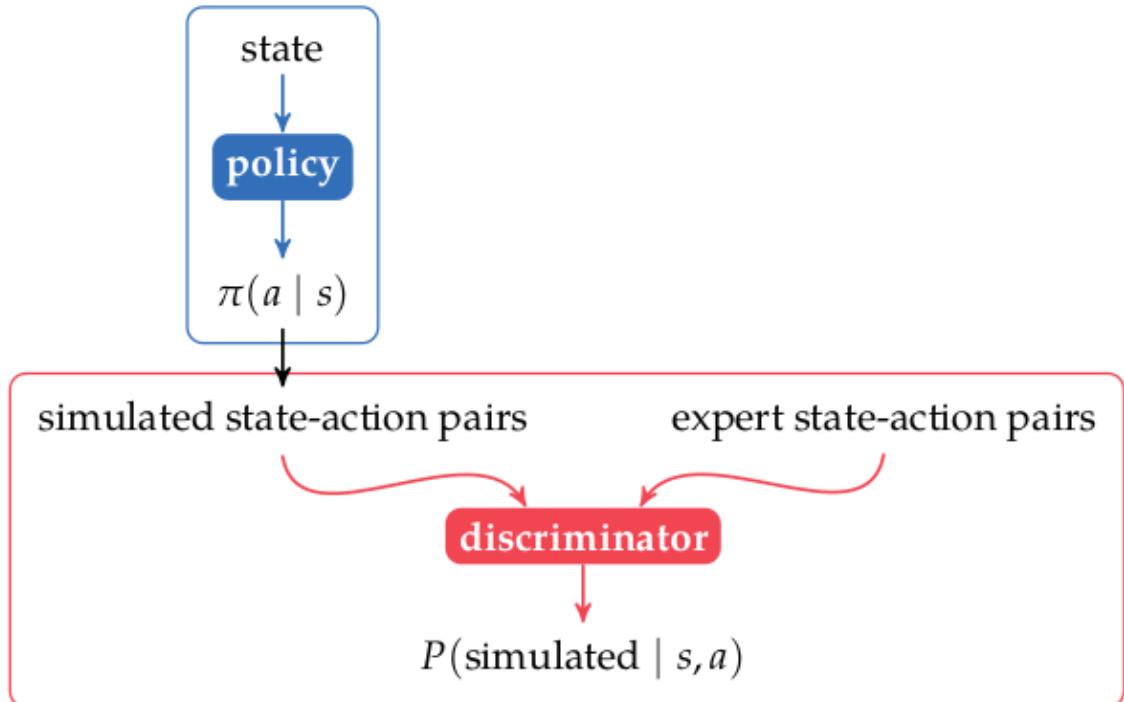


Generative Adversarial Imitation Learning (GAIL)



$$\max_{\Phi} \min_{\theta} \mathbb{E}_{(s,a) \sim \pi_\theta} [\log(C_\Phi(s, a))] + \mathbb{E}_{(s,a) \sim \mathcal{D}} [\log(1 - C_\Phi(s, a))]$$

Generative Adversarial Imitation Learning (GAIL)

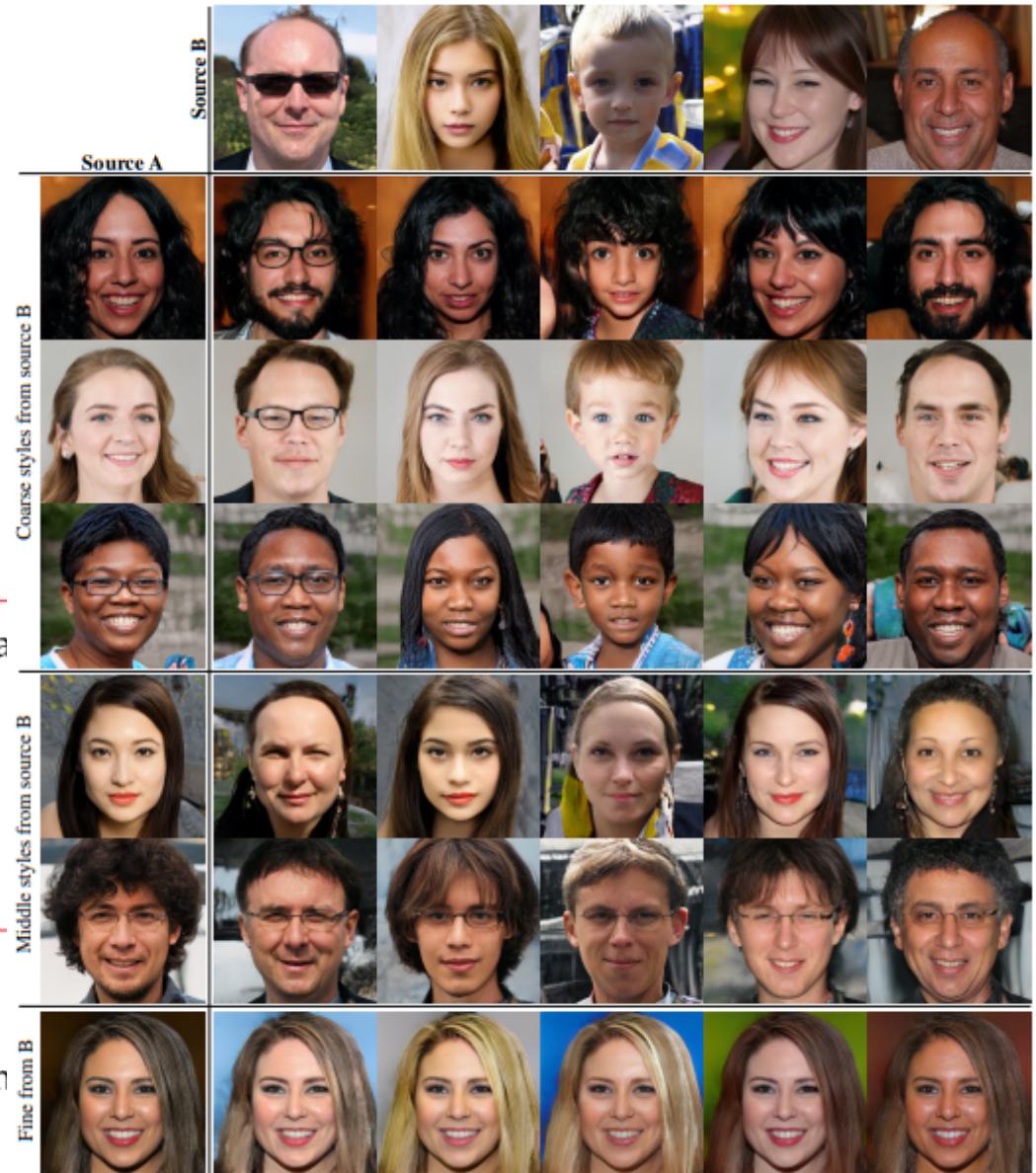


GANs are frighteningly good
at generating believable
synthetic things

$$\max_{\Phi} \min_{\theta} \mathbb{E}_{(s,a) \sim \pi_\theta} [\log(C_\Phi(s, a))] + \mathbb{E}_{(s,a) \sim \mathcal{D}} [\log(1 - C_\Phi(s, a))]$$

Adversarial Imitation Learning (GAIL)

simulation



GANs are frighteningly good
at generating believable
synthetic things

$\max_{\phi} \min_{\theta}$

Inverse Reinforcement Learning

Inverse Reinforcement Learning

What if we know the dynamics, but not the reward?

Inverse Reinforcement Learning

What if we know the dynamics, but not the reward?

Reinforcement Learning

Inverse Reinforcement Learning

Input

Output

Inverse Reinforcement Learning

What if we know the dynamics, but not the reward?

Reinforcement Learning

Input

Environment (S, A, T, R)

Inverse Reinforcement Learning

Output

π^*

$S, A, T, \{\tau\}$

R

Exercise

1	2	3
4	5	6
7	8	9

π

$$\begin{array}{ll} 1 \rightarrow & 1 \rightarrow \\ 2 \rightarrow & 2 \downarrow \\ 3 \downarrow & 5 \rightarrow \\ 6 \downarrow & 6 \downarrow \\ q & q \end{array}$$

IRL is an
underspecified
problem

What is the reward function?

Maximum Margin Inverse Reinforcement Learning

Maximum Margin Inverse Reinforcement Learning

$$R_{\Phi}(s, a) = \Phi^\top \beta(s, a)$$

Maximum Margin Inverse Reinforcement Learning

$$R_{\Phi}(s, a) = \Phi^\top \beta(s, a)$$

$$\beta(s, a) \in \{0, 1\}^n$$

Maximum Margin Inverse Reinforcement Learning

$$R_{\Phi}(s, a) = \Phi^\top \beta(s, a)$$

$$\underbrace{\beta(s, a) \in \{0, 1\}^n}_{\|\Phi\|_2 \leq 1}$$

Maximum Margin Inverse Reinforcement Learning

$$R_{\Phi}(s, a) = \Phi^\top \beta(s, a)$$

$$\beta(s, a) \in \{0, 1\}^n$$

$$\|\Phi\|_2 \leq 1$$

$$\begin{aligned}\mathbb{E}_{s \sim b}[U(s)] &= \mathbb{E}_\tau \left[\sum_{k=1}^d \gamma^{k-1} R_{\Phi}(s^{(k)}, a^{(k)}) \right] \\ &= \mathbb{E}_\tau \left[\sum_{k=1}^d \gamma^{k-1} \underbrace{\Phi^\top \beta(s^{(k)}, a^{(k)})}_{\text{Discounted expectation of feature values}} \right] \\ &= \Phi^\top \left(\mathbb{E}_\tau \left[\sum_{k=1}^d \gamma^{k-1} \beta(s^{(k)}, a^{(k)}) \right] \right) \\ &= \Phi^\top \mu_\pi\end{aligned}$$

Maximum Margin Inverse Reinforcement Learning

$$R_{\Phi}(s, a) = \Phi^\top \beta(s, a)$$

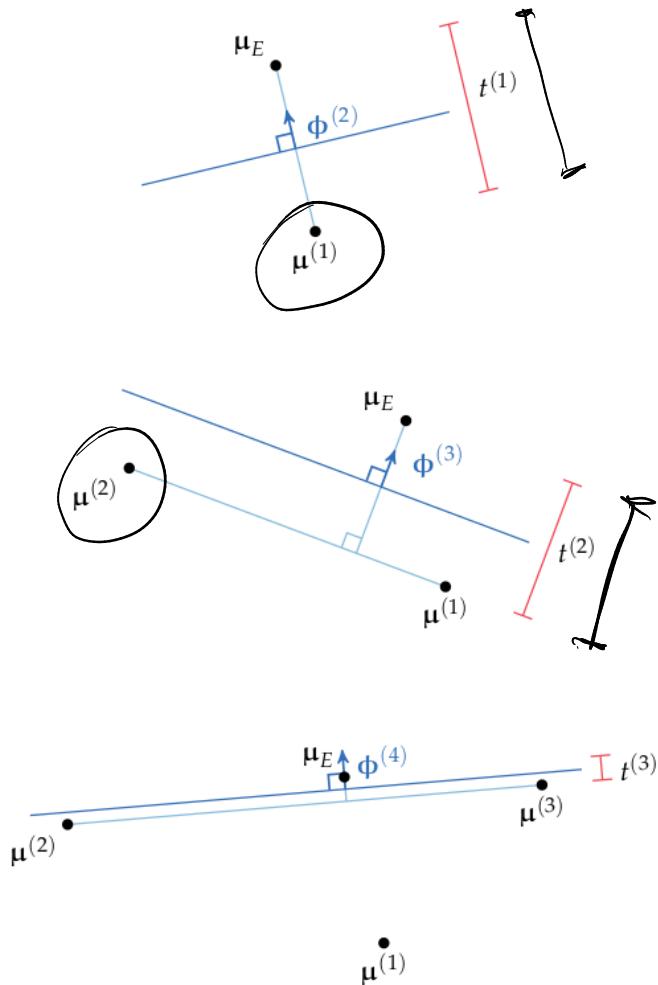
$$\beta(s, a) \in \{0, 1\}^n$$

$$\|\Phi\|_2 \leq 1$$

$$\begin{aligned} & \underset{t, \Phi}{\text{maximize}} \quad t \xrightarrow{\text{margin}} \\ & \text{subject to} \quad \underbrace{\Phi^\top \mu_E \geq \Phi^\top \mu^{(i)} + t}_{\|\Phi\|_2 \leq 1} \quad \text{for } i = 1, \dots, k-1 \end{aligned}$$

$$\begin{aligned} \mathbb{E}_{s \sim b}[U(s)] &= \mathbb{E}_\tau \left[\sum_{k=1}^d \gamma^{k-1} R_{\Phi}(s^{(k)}, a^{(k)}) \right] \\ &= \mathbb{E}_\tau \left[\sum_{k=1}^d \gamma^{k-1} \Phi^\top \beta(s^{(k)}, a^{(k)}) \right] \\ &= \Phi^\top \left(\mathbb{E}_\tau \left[\sum_{k=1}^d \gamma^{k-1} \beta(s^{(k)}, a^{(k)}) \right] \right) \\ &= \Phi^\top \mu_\pi \end{aligned}$$

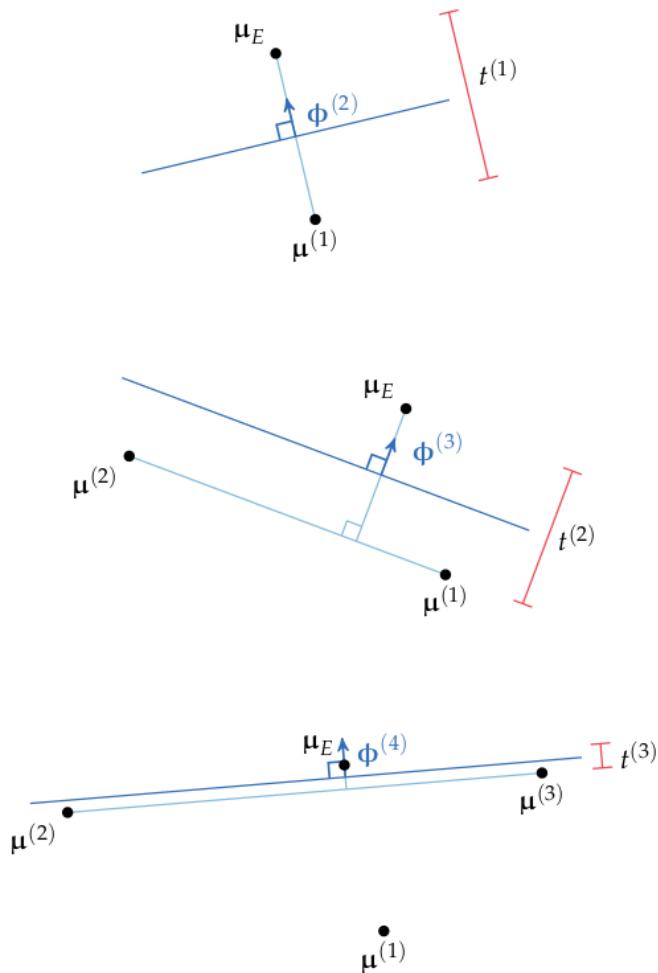
Maximum Margin Inverse Reinforcement Learning



maximize
$$t$$

subject to $\phi^\top \mu_E \geq \phi^\top \mu^{(i)} + t$ for $i = 1, \dots, k - 1$
 $\|\phi\|_2 \leq 1$

Maximum Margin Inverse Reinforcement Learning



$$\underset{t, \Phi}{\text{maximize}} \quad t$$

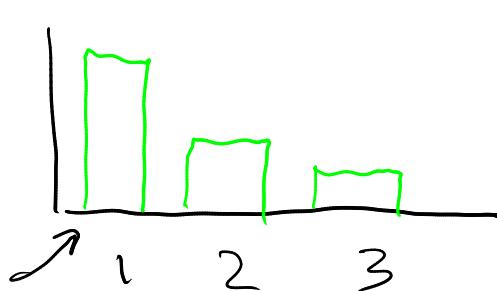
$$\text{subject to} \quad \Phi^\top \mu_E \geq \Phi^\top \mu^{(i)} + t \quad \text{for } i = 1, \dots, k-1 \\ \|\Phi\|_2 \leq 1$$

$$\underset{\lambda}{\text{minimize}} \quad \|\mu_E - \mu_\lambda\|_2$$

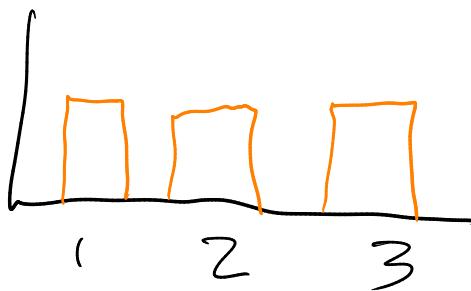
$$\text{subject to} \quad \lambda \geq 0 \\ \|\lambda\|_1 = 1$$

Principle of Maximum Entropy

$$H(X) = - \sum_x P(x) \log P(x)$$

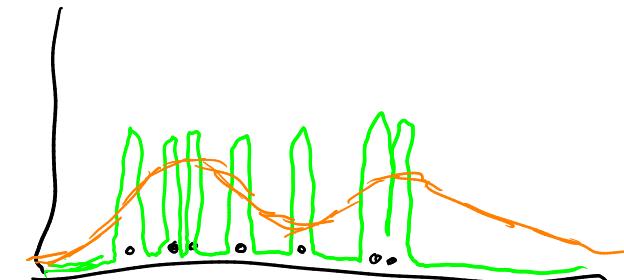
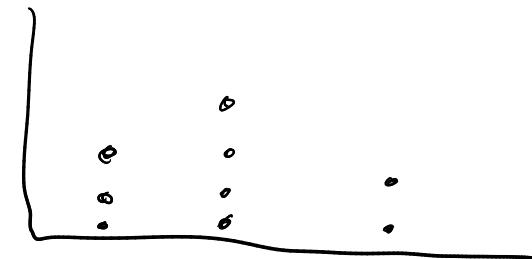
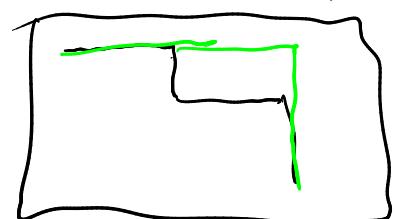


lower entropy



higher entropy

Choose R_ϕ such that $P_\phi(\tau)$
has the maximum entropy



Maximum Entropy Inverse Reinforcement Learning

Maximum Entropy Inverse Reinforcement Learning

Least informative trajectory distribution

Maximum Entropy Inverse Reinforcement Learning

Least informative trajectory distribution

$$P_{\Phi}(\tau) = \frac{1}{Z(\Phi)} \exp(\underline{R_{\Phi}(\tau)})$$

Maximum Entropy Inverse Reinforcement Learning

Least informative trajectory distribution

$$P_{\Phi}(\tau) = \frac{1}{Z(\Phi)} \exp(R_{\Phi}(\tau)) \quad Z(\Phi) = \sum_{\tau} \exp(R_{\Phi}(\tau))$$

Maximum Entropy Inverse Reinforcement Learning

Least informative trajectory distribution

$$P_{\Phi}(\tau) = \frac{1}{Z(\Phi)} \exp(R_{\Phi}(\tau)) \quad Z(\Phi) = \sum_{\tau} \exp(R_{\Phi}(\tau))$$

$$\max_{\Phi} f(\Phi) = \max_{\Phi} \sum_{\tau \in \mathcal{D}} \log P_{\Phi}(\tau)$$

Maximum Entropy Inverse Reinforcement Learning

Maximum Entropy Inverse Reinforcement Learning

$$\max_{\Phi} f(\Phi) = \max_{\Phi} \sum_{\tau \in \mathcal{D}} \log P_{\Phi}(\tau)$$

Maximum Entropy Inverse Reinforcement Learning

$$\max_{\Phi} f(\Phi) = \max_{\Phi} \sum_{\tau \in \mathcal{D}} \log P_{\Phi}(\tau)$$

Max margin IRL
 $R_{\Phi}(s,a) = \phi^T \beta(s,a)$
Max Ent R_{Φ} differentiable

$$\begin{aligned} f(\Phi) &= \sum_{\tau \in \mathcal{D}} \log \frac{1}{Z(\Phi)} \exp(R_{\Phi}(\tau)) \\ &= \left(\sum_{\tau \in \mathcal{D}} R_{\Phi}(\tau) \right) - |\mathcal{D}| \log Z(\Phi) \\ &= \left(\sum_{\tau \in \mathcal{D}} R_{\Phi}(\tau) \right) - |\mathcal{D}| \log \sum_{\tau} \exp(R_{\Phi}(\tau)) \end{aligned}$$

Maximum Entropy Inverse Reinforcement Learning

$$\max_{\Phi} f(\Phi) = \max_{\Phi} \sum_{\tau \in \mathcal{D}} \log P_{\Phi}(\tau)$$

$$\begin{aligned} f(\Phi) &= \sum_{\tau \in \mathcal{D}} \log \frac{1}{Z(\Phi)} \exp(R_{\Phi}(\tau)) \\ &= \left(\sum_{\tau \in \mathcal{D}} R_{\Phi}(\tau) \right) - |\mathcal{D}| \log Z(\Phi) \\ &= \left(\sum_{\tau \in \mathcal{D}} R_{\Phi}(\tau) \right) - |\mathcal{D}| \log \sum_{\tau} \exp(R_{\Phi}(\tau)) \end{aligned}$$

$$\nabla_{\Phi} f = \left(\sum_{\tau \in \mathcal{D}} \nabla_{\Phi} R_{\Phi}(\tau) \right) - \frac{|\mathcal{D}|}{\sum_{\tau} \exp(R_{\Phi}(\tau))} \sum_{\tau} \exp(R_{\Phi}(\tau)) \nabla_{\Phi} R_{\Phi}(\tau) \quad (18.15)$$

$$= \left(\sum_{\tau \in \mathcal{D}} \nabla_{\Phi} R_{\Phi}(\tau) \right) - |\mathcal{D}| \sum_{\tau} P_{\Phi}(\tau) \nabla_{\Phi} R_{\Phi}(\tau) \quad (18.16)$$

$$= \left(\sum_{\tau \in \mathcal{D}} \nabla_{\Phi} R_{\Phi}(\tau) \right) - |\mathcal{D}| \sum_s b_{\gamma, \Phi}(s) \underbrace{\sum_a \pi_{\Phi}(a | s)}_{\text{visitation probabilities}} \underbrace{\nabla_{\Phi} R_{\Phi}(s, a)}_{\text{optimal policy}} \quad (18.17)$$

~~choose~~

random Φ

find π_{Φ}

find $b_{\gamma, \Phi}(s)$

calc $\nabla_{\Phi} f$

update Φ with gradient descent

visitation probabilities for R_{Φ}
optimal policy

Recap

Recap

- Behavioral cloning is supervised learning to match the actions of an expert

Recap

- Behavioral cloning is supervised learning to match the actions of an expert
- A critical problem is cascading errors, which can be addressed by gathering more data with DAgger or SMILe

Recap

- Behavioral cloning is supervised learning to match the actions of an expert
- A critical problem is cascading errors, which can be addressed by gathering more data with DAgger or SMILe
- Inverse reinforcement learning is the process of learning a reward functions from trajectories in an MDP

Recap

- Behavioral cloning is supervised learning to match the actions of an expert
- A critical problem is cascading errors, which can be addressed by gathering more data with DAgger or SMILe
- Inverse reinforcement learning is the process of learning a reward functions from trajectories in an MDP
- IRL is an underspecified problem
- Maximum entropy RL solves this problem by choosing the reward function that maximizes the entropy of the trajectories of the resulting policy