

# Imitation and Inverse Reinforcement Learning

- **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?



Dean Pomerleau  
@deanpomerleau

...

1995: 2797/2849 miles (98.2%)

Replying 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 • BYTE 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.



# 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

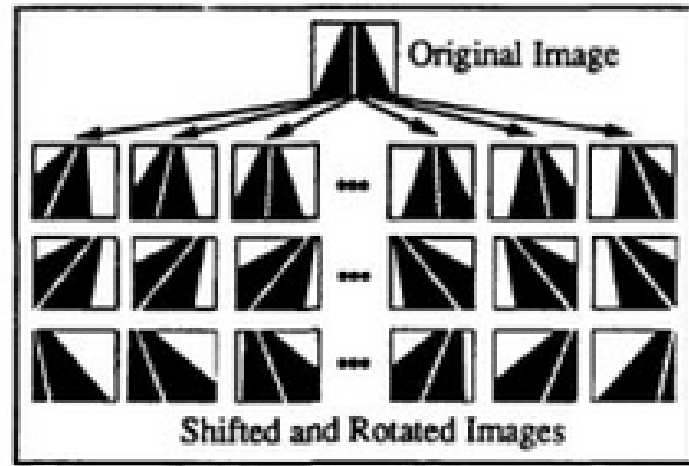


Figure 3.4: The single original video image is shifted and rotated to create multiple training exemplars in which the vehicle appears to be at different locations relative to the road.

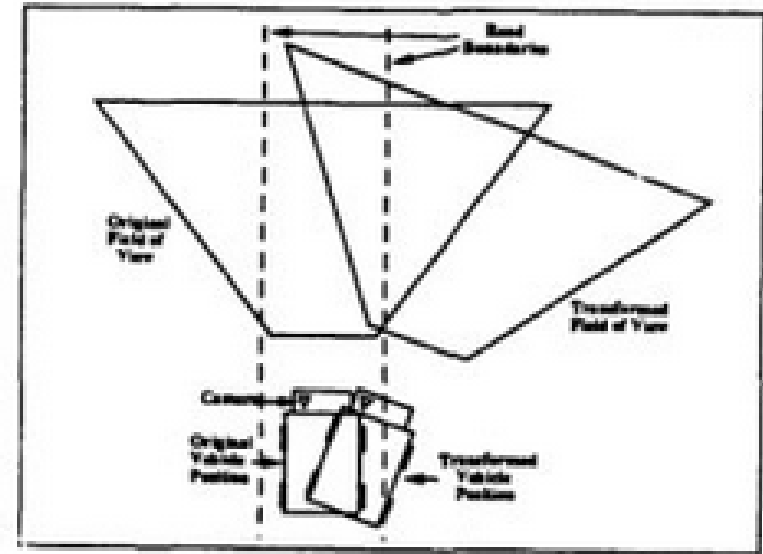


Figure 3.5: An aerial view of the vehicle at two different positions, with the corresponding sensor fields of view. To simulate the image transformation that would result from such a change in position and orientation of the vehicle, the overlap between the two field of view trapezoids is computed and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

# How did NVIDIA do it in 2016?

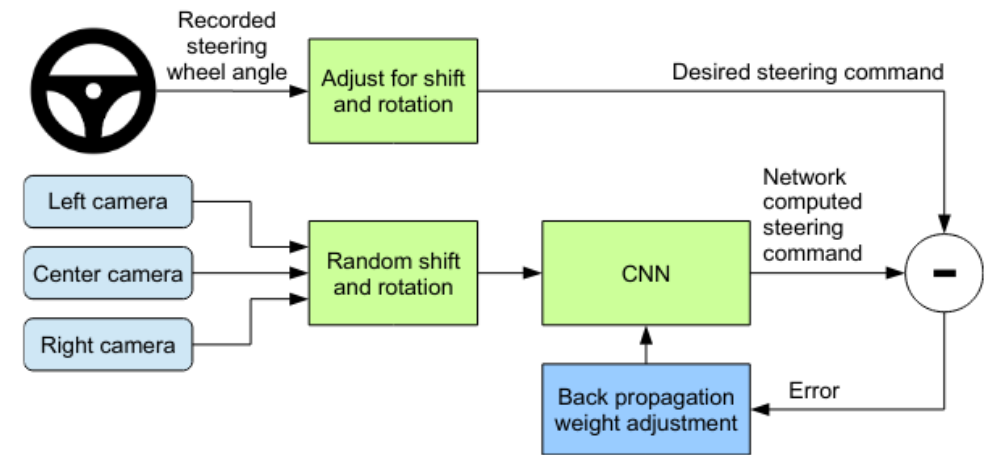
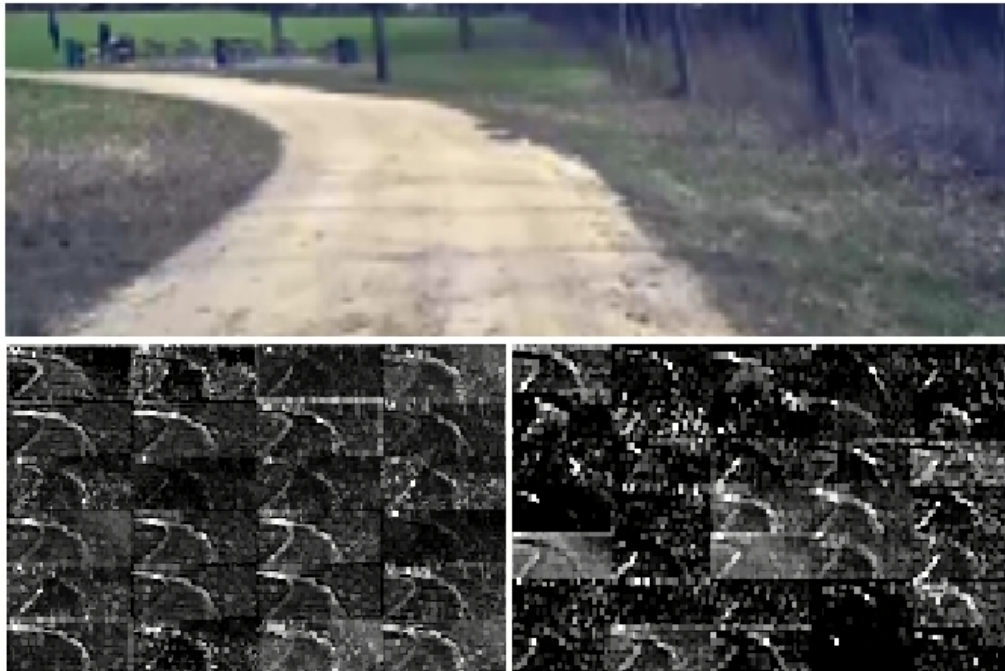


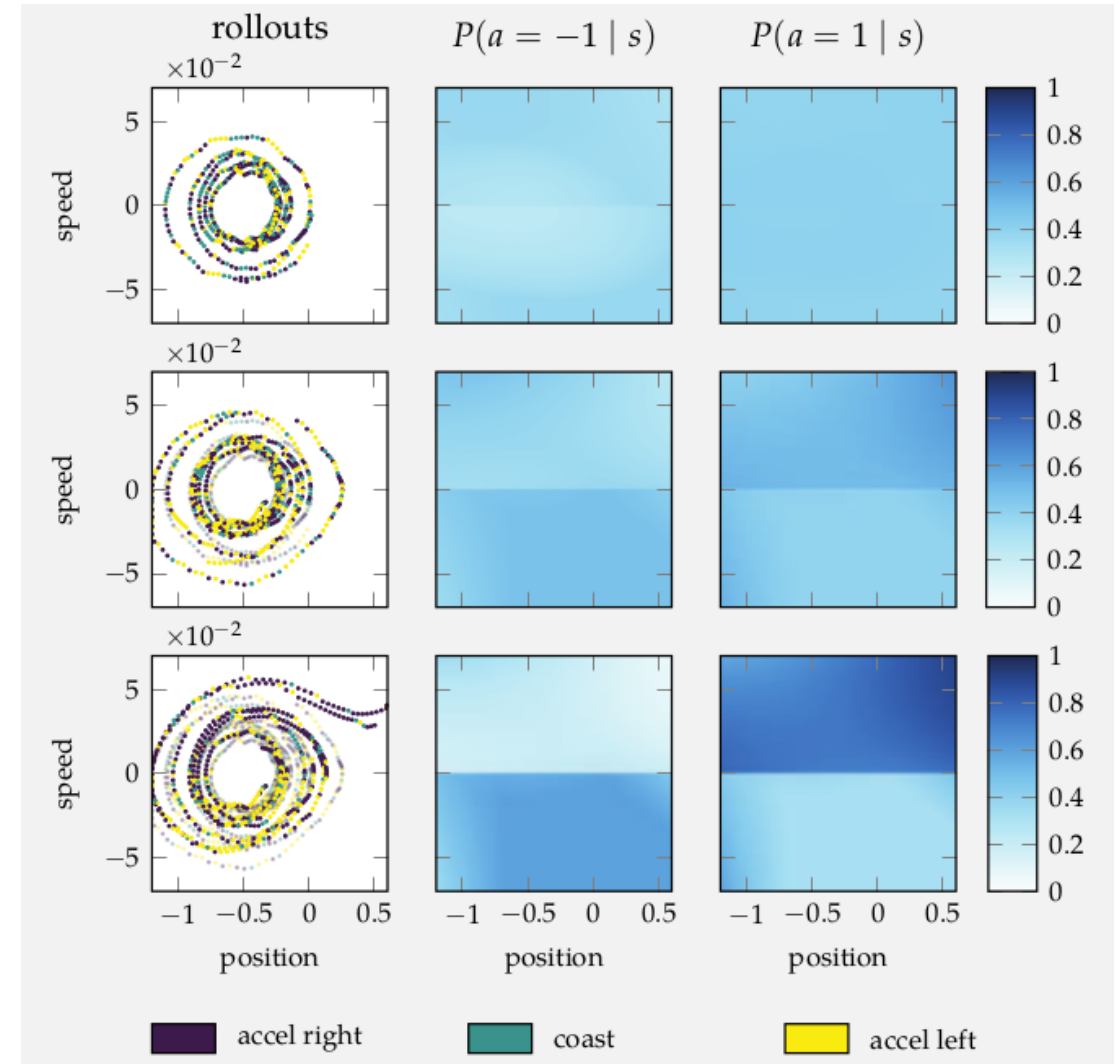
Figure 2: Training the neural network.

# Dataset Aggregation (DAgger)

```
function optimize(M::DataSetAggregation, D,  $\theta$ )  
     $\mathcal{P}$ , bc, k_max, m = M. $\mathcal{P}$ , M.bc, M.k_max, M.m  
    d, b,  $\pi_E$ ,  $\pi_\theta$  = M.d, M.b, M. $\pi_E$ , M. $\pi_\theta$   
     $\theta$  = optimize(bc, D,  $\theta$ )  
    for k in 2:k_max  
        for i in 1:m  
            s = rand(b)  
            for j in 1:d  
                push!(D, (s,  $\pi_E(s)$ ))  
                a = rand( $\pi_\theta(\theta, s)$ )  
                s = rand( $\mathcal{P}.T(s, a)$ )  
            end  
        end  
         $\theta$  = optimize(bc, D,  $\theta$ )  
    end  
    return  $\theta$   
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
```

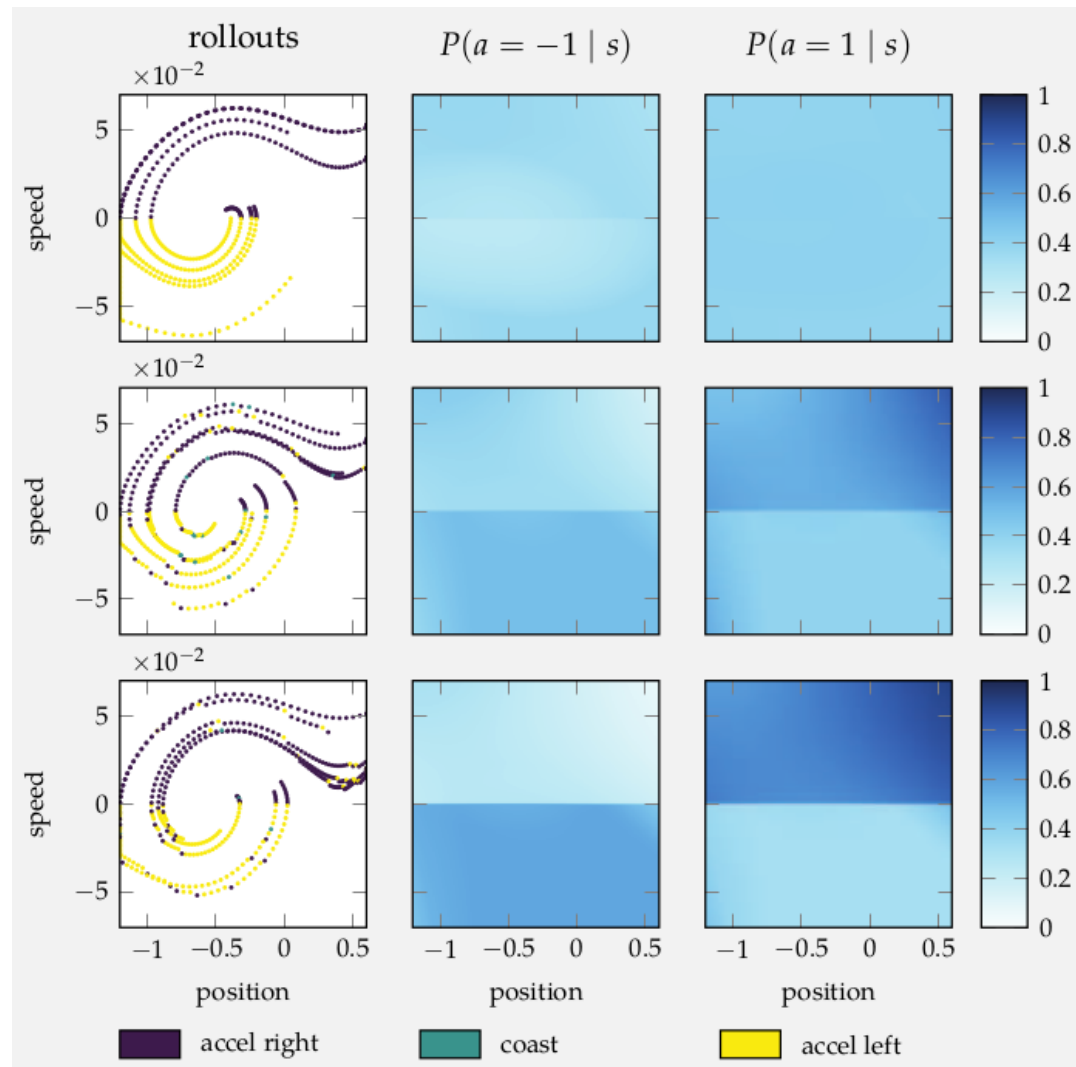
← reset D

} Gather data

← train only on D

Mix Policies

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



# Cooperative Adversarial Imitation Learning (GAIL)



GANs are frighteningly good  
at generating believable  
synthetic things



# Inverse Reinforcement Learning

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

	Reinforcement Learning	Inverse Reinforcement Learning
Input	Environment $(S, A, T, R)$	$S, A, T, \{\tau\}$
Output	$\pi^*$	$R$

# Exercise

1	2	3
4	5	6
7	8	9

$\tau$

1 →

2 →

3 ↓

6 ↓

9

1 →

2 ↓

5 →

6 ↓

9

What is the reward function?

# Principle of Maximum Entropy

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

# 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

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

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

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

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

$$= \left( \sum_{\tau \in \mathcal{D}} \nabla_{\phi} R_{\phi}(\tau) \right) - |\mathcal{D}| \sum_s b_{\gamma, \phi}(s) \sum_a \pi_{\phi}(a | s) \nabla_{\phi} R_{\phi}(s, a) \quad (18.17)$$

**Discounted visitation probability**

**Optimal policy under  $R_{\phi}$**



# 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