DSC291: Machine Learning with Few Labels

Generative Adversarial Learning

Zhiting Hu Lecture 12, May 8, 2025



HALICIOĞLU DATA SCIENCE INSTITUTE

Outline

- Deep Generative Models
 - Generative adversarial learning

- Paper presentation:
 - Letong Liang: "DeepSeek-Prover-V2"
 - Ali El Lahib, Darin Djapri: "TD-MPC2: Scalable, Robust World Models for Continuous Control"

Recap: Implicit Generative Models



Recap: Generative Adversarial Nets (GANs)

- Learning
 - A minimax game between the generator and the discriminator
 - Train D to maximize the probability of assigning the correct label to both training examples and generated samples
 - \circ Train G to fool the discriminator



Recap: Optimality of GANs

Question: in practice, we're unlikely to get the optimal D^* . In this case, what is the minimax game truly optimizing?



Theorem 1. The global minimum of the virtual training criterion C(G) is achieved if and only if $p_g = p_{data}$. At that point, C(G) achieves the value $-\log 4$.

$$C(G) = -\log(4) + KL\left(p_{\text{data}} \| \frac{p_{\text{data}} + p_g}{2}\right) + KL\left(p_g \| \frac{p_{\text{data}} + p_g}{2}\right)$$
$$= -\log(4) + 2 \cdot JSD\left(p_{\text{data}} \| p_g\right) \quad \text{Jensen-Shannon Divergence} \quad \text{Symmetric}$$
[Goodfellow et al., 2014]

Wasserstein GAN (WGAN)

If our data are on a low-dimensional manifold of a high dimensional space, the model's manifold and the true data manifold can have a negligible intersection practice

COD-Aira



Wasserstein GAN (WGAN)

- If our data are on a low-dimensional manifold of a high dimensional space, the model's manifold and the true data manifold can have a negligible intersection in practice
- The loss function and gradients may not be continuous and well behaved

Wasserstein GAN (WGAN)

- If our data are on a low-dimensional manifold of a high dimensional space, the model's manifold and the true data manifold can have a negligible intersection in practice
- The loss function and gradients may not be continuous and well behaved
- The Wasserstein Distance is well defined
 - Earth Mover's Distance
 - Minimum transportation cost for making one pile
 - of dirt in the shape of one probability distribution

to the shape of the other distribution





Progressive GAN



Progressive GAN



Progressive GAN





• GANs benefit dramatically from scaling

- GANs benefit dramatically from scaling
- 2x 4x more parameters
- 8x larger batch size
- Simple architecture changes that improve scalability

- GANs benefit dramatically from scaling
- 2x 4x more parameters
- 8x larger batch size
- Simple architecture changes that improve scalability





• GANs benefit dramatically from scaling



[Brock et al., 2018]

Key Takeaways

- Deep Generative Models: brief history
- GANs:
 - Implicit generative model \bigcirc
 - Minimax formulation $\rightarrow Jame Harry Wasserstein GAN$ Ο
 - \bigcirc



RL Conference 2024



RL Conference 2024



So far... Supervised Learning

Data: (x, y) x is data, y is label

Goal: Learn a *function* to map x -> y

Examples: Classification, regression, object detection, semantic segmentation, image captioning, etc.





Classification

So far... Unsupervised Learning

Data: x no labels!

Goal: Learn some underlying hidden *structure* of the data

Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc.



Today: Reinforcement Learning

Problems involving an **agent** interacting with an **environment**, which provides numeric **reward** signals

Goal: Learn how to take actions in order to maximize reward





Overview

- What is Reinforcement Learning?
- Markov Decision Processes
- Q-Learning
- Policy Gradients



Environment









Robot Locomotion



Objective: Make the robot move forward

State: Angle and position of the joints **Action:** Torque applied on joints **Reward:** 1 at each time step upright + forward movement

Atari Games



Objective: Complete the game with the highest score

State: Raw pixel inputs of the game stateAction: Game controls e.g. Left, Right, Up, DownReward: Score increase/decrease at each time step

Go



Objective: Win the game!

State: Position of all piecesAction: Where to put the next piece downReward: 1 if win at the end of the game, 0 otherwise



How can we mathematically formalize the RL problem?



Markov Decision Process

- Mathematical formulation of the RL problem
- Markov property: Current state completely characterises the state of the world

Defined by: $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbb{P}, \gamma)$

- ${\cal S}$: set of possible states
- ${\cal A}$: set of possible actions
- $\boldsymbol{\mathcal{R}}$: distribution of reward given (state, action) pair
- ℙ : transition probability i.e. distribution over next state given (state, action) pair
- γ : discount factor

Markov Decision Process

- At time step t=0, environment samples initial state $s_0 \sim p(s_0)$
- Then, for t=0 until done:
 - Agent selects action a_t
 - Environment samples reward $r_t \sim R(. | s_t, a_t)$
 - Environment samples next state $s_{t+1} \sim P(. | s_t, a_t)$
 - Agent receives reward r_t and next state s_{t+1}

- A policy $\pi \, \textsc{is}$ a function from S to A that specifies what action to take in each state
- **Objective**: find policy π^* that maximizes cumulative discounted reward:



A simple MDP: Grid World



Set a negative "reward" for each transition (e.g. r = -1)

Objective: reach one of terminal states (greyed out) in least number of actions

A simple MDP: Grid World





Random Policy

Optimal Policy

Questions?