nuro

Behavior Modeling and Learned Motion Selection for Safe Driving

How to incorporate Diffusion Models and Reinforcement Learning into the Autonomy Stack

Aleksandr Petiushko (apetiushko@nuro.ai) ML Research 6/19/2023

SSAD CVPR2023 Workshop

nura

The team



Content

01	Motivation
02	Diffusion: Trajectory Generation
03	RL: Motion Selection
04	Examples
05	Limitations and Conclusion

Motivation



Content

Problem 1: Road Agents

Trustworthy predictions for use in both Prediction and Simulation

Problem 2: AV Motion

Flexible and safe selection process allowing ego proposals of any source



Better training, evaluation and reasoning leading to safer driving!

Better Agents Prediction/Simulation



Usage of IL-based Prediction model for other agents can lead to unreasonable proposals due to distribution shift

Historical Approach

Use only heuristic-based agents

Better Solution

Target prediction model to **better distribution coverage/recall** (not only precision) with some ways to use it for getting good minADE

Long-Horizon Planning by Selection

Problem

Decisions have long-term, delayed consequences

Historical Approach

Use long-term **predictions** to approximate long-horizon planning

Better Solution

Learn a model that takes into account an **expectation over all futures**. Selection to narrow down the search space Diffusion models: Background and Trajectory Generation



8

What is a Diffusion Model?

Diffusion Model: it is a generative model (markovian hierarchical variational autoencoder)

- Adding step by step some portion of noise as a diffusion analogy
- Forward diffusion process: adding noise by $q(x_t|x_{t-1})$. Also known as *encoding*
- **Reverse** diffusion process: de-noising by $p(x_{t-1}|x_t)$. Also known as *decoding*



Image credit: https://arxiv.org/pdf/2208.11970.pdf

Success of Diffusion Models



Text2Image (along with audio, video) generation: Done! (sorry, GANs 😢)

But what about other tasks?

Diffusion Models for Autonomous Driving

But what about other tasks?



We are combining both functionalities: prediction and simulation

Zhong, Ziyuan, et al. "<u>Guided Conditional Diffusion for Controllable Traffic Simulation</u>", 2022 Jiang, Chiyu, et al. "<u>MotionDiffuser: Controllable Multi-Agent Motion Prediction using Diffusion</u>", 2023 (01)

DTG: Main Goals

Development of Trajectory Generation module capable of a good distribution coverage

DTG = Diffusion-based Trajectory Generator



Improvement of closed-loop simulations

DTG: Main Goals



Development of Trajectory Generation module (decoder) capable of a good distribution coverage Is theoretically ensured by using Variational Diffusion Model (VDM) by explicit ELBO (~NLL) optimization

DTG: Main Goals

Will provide more useful signal for RL-based trainings

(02)

Improvement of closed-loop simulations

DTG: Features



Learn diverse behaviors with distribution that matches real-world driver behaviors (02)

Provide good NLL, minADE, and other Prediction-aware metrics (03)

Lead to stable, consistent and realistic simulation

VDM for Trajectory Generation

DTG: Current Architecture



Note: we can use different encoders (lstm-based, transformer-based)

DTG: Ensuring good minADE

- Vanilla **VDM** models the distribution of trajectories, the sampled *N* trajectories not necessarily have 1 close to GT
- We can mitigate it through clustering for getting a good minADE
 - And even probability as a size of cluster!



Kingma, Diederik, et al. "Variational diffusion models", 2021

RL: Background and Motion Selection



What is Deep RL?



What is Deep RL?



What is Deep RL?



How to Optimize?

Objective: maximize reward under the policy while limiting probability of risky events Learn: state-action Q value function

Optimize: iteratively improve *Q* for all *s* and *a*

What does this look like?

0.43	0.48	0.53	0.59	0.66	0.73	— 0	0	0	0	1	l
0.39	0.43	0.48	0	0.73	0.81	•	0	•	0	0.59	0.53
0.35	0.39	0.66	0.73	0.81	0.9	1	0.9	0.81	0.73	0.66	0.59
0	0.66	0.73	0.81	0.9	1	0 R:1	ļ	0.9	0.81	0.73	0.66
0.53	0.59	0.29	0.26	0.81	0.9	1	0.9	0.81	0.73	0.66	0.59
0.48	0.53	0.26	0.66	0.73	0.81	0.9	0.81	0.73	0.66	0.59	0.53
0.43	0.48	0.43	0.59	0.66	0.73	0.81	0.73	0.66	0.59	0.53	0.48

Image credit: https://towardsdatascience.com/interactive-g-learning-9d9203fdad70

RL for Selection

Why Motion Selection?

(01)

Discrete problem. Rank trajectories rather than produce them.

RLMS = RL for Motion Selection (02)

Low-level decision making well handled by trajectory generation modules

Why Motion Selection?



Discrete problem. Rank trajectories rather than produce them. Allow heuristics and domain knowledge to filter the trajectory space for RL (02)

Low-level decision making well handled by trajectory generation modules

Anatomy of the RLMS Model: Basic RL



Basic RL: Limitations



No concrete notion or **constraint** on safety

Anatomy of the RLMS Model: Risk Sensitive RL



Risk Sensitive RL: Limitations

(01)

No hard constraint on safety

Anatomy of the RLMS Model: Constrained RL



Anatomy of the RLMS Model: Constrained RL





RLMS: Block Diagram

First execute current RLMS policy in the simulator and store trajectories

RLMS Policy Sim 1 Sim 2 Sim 3 Replay Buffer

RLMS: Block Diagram

Use saved trajectories to train task critic and risk critic

Replay RLMS Sim 3 Sim 1 Sim 2 Policy Buffer Reward 1 Reward 2 Task Critic Reward 3 Transitions Risk 1 Risk Risk 2 Critic Risk 3

RLMS: Block Diagram

Combine task and risk critic values into utilities and train policies for each

Replay RLMS Sim 1 Sim 2 Sim 3 Buffer Policy Reward 1 Reward 2 Task Task Task Util Critic Policy Reward 3 Transitions Risk 1 Recovery Recovery Risk 2 Risk Util Policy Critic Risk 3

Combine task and recovery policy to get RLMS policy

RLMS: Block Diagram



Constructing RLMS Mixed Policy with Recovery RL

The final RLMS policy uses a combination of both the task and recovery policies to sample actions.

We first score all possible actions with each of our risk critics



Constructing RLMS Mixed Policy with Recovery RL

The final RLMS policy uses a combination of both the task and recovery policies to sample actions.

If there exist safe actions then sample from re-normalized task policy



Constructing RLMS Mixed Policy with Recovery RL

The final RLMS policy uses a combination of both the task and recovery policies to sample actions.

Otherwise sample from recovery policy



Anatomy of the RLMS Model



Examples





OK: Vehicle overtaking NuroBot on the left



Top: Onroad log

Bottom: Sim

Video link: https://www.youtube.com/watch?v=FE7IR11uVB8

OK: Occluded Unprotected Left



Top: Onroad log

Bottom: Sim

Video link: https://www.youtube.com/watch?v=scblFi50oA8

Not OK: Problems with stability - selection is a combination of plans because we don't have a single initial good source to choose from (making up its own plan via flicker yield)



Top: Onroad log

Bottom: Sim

Video link: https://www.youtube.com/watch?v=FE7IR11uVB8

Limitations and conclusions



DTG: Limitations



Sampling-only inference (hard to use in the production) (02)

Latency-performance tradeoff

(03)

Non-deterministic simulation

RLMS: Limitations



Still no *hard* constraint on safety



Rare sparse events still challenging to learn (i.e. collisions)

(03)

Sample inefficient – takes many simulation steps to learn

Conclusions



Diffusion-based models help to match the distributions, not points



Learning selection provides long-horizon reasoning (03)

Recent academic SotA can be used for practical tasks to add more safety!

Join ML Research!

Open Roles			
Machine Learning Research	~	Current Roles	
Mountain View, California (HQ)	~	Software	^
Full-Time	~	Machine Learning Research	
Search	Q	Machine Learning Research Scientist Mountain View California (HO)	Full-Time
Clear Filters		Machine Learning Research Scientist - Reinforcement Learning Mountain View, California (HQ)	Full-Time

https://www.nuro.ai/careers