

# Motion Planning for Autonomous Driving

Philippe Weingertner, Minnie Ho

November 8, 2019

## 1 Introduction

## 2 Related Work

There is a rich literature related to Motion Planning and a very detailed survey of traditional methods is provided in [14]. Among the first 4 successful participants of DARPA Urban Challenge in 2007, the approaches vary. The winner, CMU Boss vehicle used variational techniques for local trajectory generation in a structured environment. This was done in a 2 steps path-velocity decomposition. A first step of path planning, using variational techniques, is performed and for every candidate path, a combination of different velocity profiles (constant, linear, linear ramp, trapezoidal) is applied : the combination of a path and velocity profile defines a trajectory. In unstructured environments (parking lots) or in error recovery situations a lattice graph in 4-dimensional configuration space (position, orientation and velocity) is searched with Anytime D\* algorithm to find a collision-free path. More details are provided in [6, 1, 12]. The vehicle from Stanford used a search strategy coined Hybrid A\* that constructs a tree of motion primitives by recursively applying a finite set of maneuvers. The search was guided by a carefully designed heuristic. The vehicle arriving 3rd, Victor Tango from Virginia Tech, constructs a graph discretization of possible maneuvers and searches the graph with the A\* algorithm. The vehicle arriving 4th, developed by MIT used a variant of RRT algorithm with biased sampling. While all these techniques differ, they fundamentally rely on a graph search where nodes correspond to a configuration state and edges correspond to elementary

motion primitives. Although they provide solutions, the runtime and state space can grow exponentially large. In this context, the use of heuristic to guide the search is important. **(NB: uncertainty is not properly considered here in these traditional MP methods. It is like reasoning with the mean state vector provided by sensor fusion output and ignoring the covariance matrix !)**

More recently, Reinforcement Learning and Deep RL have been investigated in the context of Autonomous Driving for Decision Making either at the Behavioural Planning or Motion Planning level. In some research papers from Volvo [8] and BMW [7], an RL agent is trained in simulation to take decision at a higher tactical level: the decisions relate to a maneuver selection, like lane change, rather than a low level acceleration command. DQN is used to train an agent. But the problem with Reinforcement Learning is that the utility is optimized in expectation. So even if the reward is designed to avoid collisions, this will be optimized in expectation: ultimately it is as if safety would be enforced with soft constraints rather than hard constraints. Which is of course not acceptable for a real vehicle. To solve this problem in [7] an additional safety check layer is added after the DQN agent to eventually override the DQN agent decision if it is considered unsafe. Checking a decision wrt to a specific criteria is simpler than designing a decision making system that jointly optimizes efficiency, comfort and safety objectives. With RL applied to AD we have to account for additional safety checks. In [2] Deep RL is applied at the local planner level: the action space is a set of longitudinal accelerations  $\{-4m/s^2, -2m/s^2, 0m/s^2, 2m/s^2\}$  applied along a given path at a T-intersection. Safety is

handled in a different way here compared to previous BMW approach: the agent is constrained to choose among a restricted set of safe actions per state. So the safety is enforced before Deep RL. Ultimately car manufacturers may want to combine both types of safety checks: constraining the action set per state before enabling an RL agent to make its own decision, and checking again the final sequence of decisions proposed by the RL agent.

Now the interesting topic is how to best combine traditional Motion Planning with RL. What are the limitations of these techniques in isolation and how to use the strengths of both approaches and circumvent their weaknesses. Traditional motion planning relies heavily on tree search and to enable real time solutions good heuristics are required. Designing a good heuristic is hard. What if we could learn it? By training an agent with model free RL we can potentially end up with an agent that performs pretty well most of the time and from time to times fails miserably in a way that is hard to explain. The main problems with model free RL are sample efficiency (we need a lot of data), enforcing hard constraints and explainability (how can we explain the decision taken by a RL agent which may become a problem for a car manufacturer). While a model based planning method has the advantages of explainability, do not rely on data and can deal in a more systematic way with hard constraints. As demonstrated in [3] in simple situations RL methods have no benefit over rule based methods, pure RL does not enable the agent to act in a safer way. But when the situation becomes much more complex with an increasing number of cars and pedestrians, the benefits of Deep RL methods become clear.

In the gaming domain, chess and go, performances superior to human performances have been achieved with AlphaGo Zero [18]: by combining planning with MCTS tree search and learning with RL. A neural network biases the sampling towards the most relevant parts of the search tree: a learnt policy-value function is used as a heuristic during inference. While during training, MCTS is used to improve the sample efficiency of RL training. Now there are a few major differences between a game like chess or go and our initial Motion Planning problem. In chess or go

the state space is discrete and fully observable while in AD the state space is continuous and partially observable. In terms of action sets in both cases, we can deal with discrete action sets. But another challenge is that self-play can not be used in the context of Motion Planning. These challenges have been recently tackled in different publications. The applicability of AlphaGo Zero to Autonomous Driving has been studied in [9, 4, 15].

TODO:

- analysis of the 3 most relevant papers
- highlight what is specific to our use case

### 3 Approach

A few notes + baseline / oracle recap + explain our Custom openai-gym env

#### 3.1 MDP model

WARNING: use a relative representation of the state (to reduce state space “area”, enable easier generalization/learning etc ...)

- **State:**  $S_t = \{S_i^t\}_{i=0..11} = \left\{ (x, y, v_x, v_y)_{\text{ego}}, (x, y, v_x, v_y)_{\text{obj}_{1..10}} \right\}$   
 – with  $S_i^t = [x, y, v_x, v_y]^T$  and  $i \in [0, 11]$
- **Actions:**  $a \in [-2 \text{ m s}^{-2}, -1 \text{ m s}^{-2}, 0 \text{ m s}^{-2}, 1 \text{ m s}^{-2}, 2 \text{ m s}^{-2}]$   
 – for ego vehicle we choose an acceleration along y-axis
- **Transitions:**  $T(s' | s, a) = P(S_i^{t+1} | S_i^t, a) = \mathcal{N}\left(T_s S_i^t + T_a \begin{bmatrix} a_x \\ a_y \end{bmatrix}\right)$   
 – Linear Gaussian dynamics with a Constant Velocity Model  
 –  $T_s = \begin{bmatrix} 1 & 0 & \text{dt} & 0 \\ 0 & 1 & 0 & \text{dt} \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}, T_a = \begin{bmatrix} \frac{\text{dt}^2}{2} & 0 \\ 0 & \frac{\text{dt}^2}{2} \\ \text{dt} & 0 \\ 0 & \text{dt} \end{bmatrix}$

$$- \begin{bmatrix} x^{t+1} \\ y^{t+1} \\ v_x^{t+1} \\ v_y^{t+1} \end{bmatrix} = \begin{bmatrix} x^t + v_x^t dt \\ y^t + v_y^t dt \\ v_x^t \\ v_y^t \end{bmatrix} + \begin{bmatrix} a_x \frac{dt^2}{2} \\ a_y \frac{dt^2}{2} \\ a_x dt \\ a_y dt \end{bmatrix} + \begin{bmatrix} \sigma_x \\ \sigma_y \\ \sigma_{v_x} \\ \sigma_{v_y} \end{bmatrix}$$

- **Reward:** efficiency + safety + comfort

$$- R_t = -1 - 1000 \times 1 [\text{d(ego,obj)} \leq 10] - 1 [|a_t| = 2]$$

### 3.2 Algo 1, MCTS tree search

The MDP is solved online via MCTS tree search. This planning method, RL Model based, has a complexity that does not grow exponentially with the horizon. TODO: our own implementation.

```

1: function SELECTACTION( $s, d$ )
2:   loop
3:     SIMULATE( $s, d, \pi_0$ )
4:   end loop
5:   return  $\arg \max_a Q(s, a)$ 
6: end function
1: function SIMULATE( $s, d, \pi_0$ )
2:   if  $d = 0$  then
3:     return 0
4:   end if
5:   if  $s \notin T$  then
6:     for  $a \in A(s)$  do
7:        $(N(s, a), Q(s, a)) \leftarrow (N_0(s, a), Q_0(s, a))$ 
8:     end for
9:      $T = T \cup \{s\}$ 
10:    return ROLLOUT( $s, d, \pi_0$ )
11:   end if
12:    $a \leftarrow \arg \max_a Q(s, a) + c \sqrt{\frac{\log N(s)}{N(s, a)}}$ 
13:    $(s', r) \sim G(s, a)$ 
14:    $q \leftarrow r + \lambda \text{SIMULATE}(s, d - 1, \pi_0)$ 
15:    $N(s, a) \leftarrow N(s, a) + 1$ 
16:    $Q(s, a) \leftarrow Q(s, a) + \frac{q - Q(s, a)}{N(s, a)}$ 
17:   return  $q$ 
18: end function
1: function ROLLOUT( $s, d, \pi_0$ )
2:   if  $d = 0$  then
3:     return 0
4:   end if
5:    $a \sim \pi_0(s)$ 
6:    $(s', r) \sim G(s, a)$ 
7:   return  $r + \lambda \text{ROLLOUT}(s', d - 1, \pi_0)$ 
8: end function

```

TODO: add progressive widening to deal with a continuous state space.

### 3.3 Algo 2, Approximate Q-learning

Model Free RL method

- $\hat{Q}_{\text{opt}}(s, a; \mathbf{w}) = \mathbf{w} \cdot \phi(s, a)$

- $\hat{V}_{\text{opt}}(s') = \max_{a' \in \text{Actions}(s')} \hat{Q}_{\text{opt}}(s', a')$

- Objective =  $\left( \hat{Q}_{\text{opt}}(s, a; \mathbf{w})_{\text{pred}} - \left( r + \gamma \hat{V}_{\text{opt}}(s') \right)_{\text{targ}} \right)^2$

- $\mathbf{w} \leftarrow \mathbf{w} - \eta \left[ \hat{Q}_{\text{opt}}(s, a; \mathbf{w})_{\text{pred}} - \left( r + \gamma \hat{V}_{\text{opt}}(s') \right)_{\text{targ}} \right] \phi(s, a)$

### 3.4 Algo 3, Deep Q-learning

Model Free RL method. TODO post progress report.

### 3.5 Algo 4, MCTS tree search with a learned heuristic

Combining Planning and Learning, Model Based and Model Free RL. TODO post progress report.

## 4 Experimental Setup and Status

The source code is available here: CS221 Project

- Baseline: simple rule - reflex based
- Oracle: assumes no uncertainty, UCS/A\* tree search
- Sequential Decision Making with Uncertainty => solve a MDP
  - Planning with MCTS tree search
  - Approximate Q-learning
  - Deep Q-learning
  - Combining Planning and Learning (with an efficient learned heuristic)

## References

- [1] David Bissell et al. “Autonomous automobiles: The social impacts of driverless vehicles”. In: *Current Sociology* (Dec. 2018), p. 001139211881674. DOI: 10.1177/0011392118816743.

- [2] Maxime Bouton et al. “Reinforcement Learning with Probabilistic Guarantees for Autonomous Driving”. In: *CoRR* abs/1904.07189 (2019). arXiv: 1904.07189. URL: <http://arxiv.org/abs/1904.07189>.
- [3] Maxime Bouton et al. “Safe Reinforcement Learning with Scene Decomposition for Navigating Complex Urban Environments”. In: *CoRR* abs/1904.11483 (2019). arXiv: 1904.11483. URL: <http://arxiv.org/abs/1904.11483>.
- [4] Panpan Cai et al. “LeTS-Drive: Driving in a Crowd by Learning from Tree Search”. In: *CoRR* abs/1905.12197 (2019). arXiv: 1905.12197. URL: <http://arxiv.org/abs/1905.12197>.
- [5] Michael Everett, Yu Fan Chen, and Jonathan P. How. “Motion Planning Among Dynamic, Decision-Making Agents with Deep Reinforcement Learning”. In: *CoRR* abs/1805.01956 (2018). arXiv: 1805.01956. URL: <http://arxiv.org/abs/1805.01956>.
- [6] Dave Ferguson, Thomas M. Howard, and Maxim Likhachev. “Motion Planning in Urban Environments”. In: *The DARPA Urban Challenge*. 2009.
- [7] Andreas Folkers, Matthias Rick, and Christof BäEskens. “Controlling an Autonomous Vehicle with Deep Reinforcement Learning”. In: June 2019. DOI: 10.1109/IVS.2019.8814124.
- [8] Carl-Johan Hoel, Krister Wolff, and Leo Laine. “Automated Speed and Lane Change Decision Making using Deep Reinforcement Learning”. In: *CoRR* abs/1803.10056 (2018). arXiv: 1803.10056. URL: <http://arxiv.org/abs/1803.10056>.
- [9] Carl-Johan Hoel et al. “Combining Planning and Deep Reinforcement Learning in Tactical Decision Making for Autonomous Driving”. In: *CoRR* abs/1905.02680 (2019). arXiv: 1905.02680. URL: <http://arxiv.org/abs/1905.02680>.
- [10] Mykel J. Kochenderfer. *Decision Making Under Uncertainty: Theory and Application*. MIT Press, 2015.
- [11] Edouard Leurent. “A Survey of State-Action Representations for Autonomous Driving”. working paper or preprint. Oct. 2018. URL: <https://hal.archives-ouvertes.fr/hal-01908175>.
- [12] M. McNaughton et al. “Motion planning for autonomous driving with a conformal spatiotemporal lattice”. In: *2011 IEEE International Conference on Robotics and Automation*. May 2011, pp. 4889–4895. DOI: 10.1109/ICRA.2011.5980223.
- [13] Subramanya Nagesh Rao, H. Eric Tseng, and Dimitar P. Filev. “Autonomous Highway Driving using Deep Reinforcement Learning”. In: *CoRR* abs/1904.00035 (2019). arXiv: 1904.00035. URL: <http://arxiv.org/abs/1904.00035>.
- [14] B. Paden et al. “A Survey of Motion Planning and Control Techniques for Self-Driving Urban Vehicles”. In: *IEEE Transactions on Intelligent Vehicles* 1.1 (Mar. 2016), pp. 33–55. DOI: 10.1109/TIV.2016.2578706.
- [15] F. Pusse and M. Klusch. “Hybrid Online POMDP Planning and Deep Reinforcement Learning for Safer Self-Driving Cars”. In: *2019 IEEE Intelligent Vehicles Symposium (IV)*. June 2019, pp. 1013–1020. DOI: 10.1109/IVS.2019.8814125.
- [16] Markus Schratter et al. “Pedestrian Collision Avoidance System for Scenarios with Occlusions”. In: *CoRR* abs/1904.11566 (2019). arXiv: 1904.11566. URL: <http://arxiv.org/abs/1904.11566>.
- [17] David Silver et al. “Mastering Chess and Shogi by Self-Play with a General Reinforcement Learning Algorithm”. In: *CoRR* abs/1712.01815 (2017). arXiv: 1712.01815. URL: <http://arxiv.org/abs/1712.01815>.
- [18] David Silver et al. “Mastering the game of Go without human knowledge”. In: *Nature* 550 (2017), pp. 354–359.