

---

# Reinforcement Learning with Hard Constraints for Autonomous Driving

---

**Philippe Weingertner, Vaishali G Kulkarni**  
pweinger@stanford.edu, vaishali@stanford.edu

**Project Mentor: Ramtin Keramati**

## 1 Questions

### 1.1 What is the problem?

**Why is it important? What are the key limitations of prior work? What are you proposing to do? What will that allow us to do now? (e.g. it addresses some prior key limitations, is computationally faster, etc etc)**

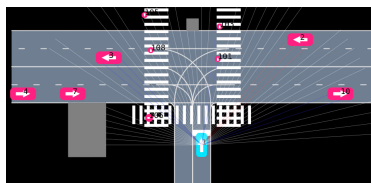
We consider the problem of decision making for an autonomous car. The Autonomous Driving pipeline consists of 3 parts:

1. Perception: sensors, localization and sensors fusion provide a world model that will be used for scene understanding and decision making.
2. Planning: based on the above world model, decisions like accelerate, slow down or change lane are taken and a trajectory is planned.
3. Control: the planned trajectory is followed as close as possible taking into account a detailed dynamical vehicle model and actuators command are sent to control the vehicle.



We focus on the Planning module. This module is typically further decomposed into 3 sub-modules:

1. A prediction module: predicting the trajectories and intents of other drivers. Typically a Driving Model can be inferred in real time to match the behavior of other drivers and to anticipate what they could do in the future.
2. A decision making module: the decisions are typically abstract and higher level like change lane but can also be low level like change longitudinal or lateral accelerations.
3. A motion planning module: based on above higher level decision, a trajectory will be planned



The difficulty in Decision Making for Autonomous Driving is due to the combination of different challenges but the main challenge we would like to address here is how to make good decisions in a situation for which we have:

- No established driving model for other vehicles. Every driver is different and may behave differently. So by construction we are in a model free setting. We do not know the driving model of other drivers. We can try to predict such a model, to estimate it, but in reality there is no absolute ground truth. And the way someone is driving may change suddenly. This is the main difficulty: we are faced with a **model free** setting, with **uncertainties** and **non stationary behaviors**.
- **Safety requirements**: we have to make safe decisions to avoid collisions and not just 99% of the time. Nevertheless there are 2 type of collisions, responsible and non-responsible ones. There will never be a guarantee for 0% collisions, but we have to avoid “responsible collisions” and in the remaining cases, for non-responsible collisions, we have to make decisions that diminish the consequences of a collision.

This problem is of paramount importance and very challenging. The model free setting, where we do not know the driving model of others, is a good match for a RL model free formulation of the problem. But the safety requirements are handled with rules we want to enforce and check. With a RL solution we are defining an objective that is optimized in expectation. It can be very good in expectation. But if it has variance, we may occasionally fail the safety requirements. So in terms of problem formulation we would like to add hard constraints to an objective in expectation.

Now let’s consider further what these Safety Constraints could be. Ultimately we want to avoid collisions. But this is a very late signal and we want to use a signal, that we can constrain and use ahead of time to prevent collisions or to diminish the consequences of collisions (in case they can not be avoided). A signal of interest is the Time To Collision: it is a sort of a proxy signal for the collisions we can explicit. Once available we can define constraints on it. So typically the constraint we will use is that the minimum Time To Collision shall be above some margin. As we are dealing with uncertainties, we do not know the driving models of other, the TTC or minTTC is actually a Random Variable. So ultimately in the way we compute the TTC or minTTC and the way we set constraints, we have to account for uncertainties. As a quick summary:

1. An AD pipeline typically predicts the Driver Model or trajectories of others. Let’s assume it tries to best match an IDM/MOBIL driver model for every surrounding car. The IDM driver model enables to predict longitudinal acceleration, it depends on 5 parameters, whereas the MOBIL driver model enables to predict the lateral acceleration and depends on 3 parameters (one parameter being a level of politeness). These models are much more valuable than raw trajectories prediction, because as input they take into account contextual information: like relative speeds or distances. So they take into account what other cars do, to predict what the car of interest will do. And if the context change, the prediction will change. Whereas simple models like CV (Constant Velocity) or CA (Constant Acceleration) do not take into account contextual information. So typically CV or CA models may be reliable for very short time horizons whereas IDM/MOBIL models are very usefull for longer time horizons or when a driver is doing something more complex like a lane change or lane merge.
2. Once we have driving models available for surrounding cars, we can evaluate the minTTC. To simplify things, assume we are in a scene with 2 cars driving longitudinally. The prediction module has best fitted 2 IDM models with 5 parameters each. And estimated the range of uncertainties for every parameter. Based on this we can define a conservative or robust TTC value: we sample possible trajectories and account for the min TTC obtained.

And ultimately we want to have a constraint on this minTTC: we would like to enforce that  $\minTTC \geq \text{margin}$

Now this is where the main challenge is when considering RL with Hard Constraints. The constraint that is of interest to us here, is very complex to compute. There is no simple differentiable graph that could be defined to express this constraint. With above IDM models: we would iterate over time steps for every cars, and once we know the new position of every other car we can estimate the longitudinal acceleration of one car and so on ... And on top of that to account for uncertainties we sample IDM parameters. So the hard constraint computation is a non trivial piece of code that is not well suited for expression as a Tensorflow Graph. Most of the papers dealing with RL and hard constraints assume the constraint is expressed as a mathematical expression that is differentiable. So all the methods dealing with Lagrangian formulation or Lyapunov stability analysis are not applicable to our case. There tend to be lots of mathematical derivations in most of the papers, but ultimately

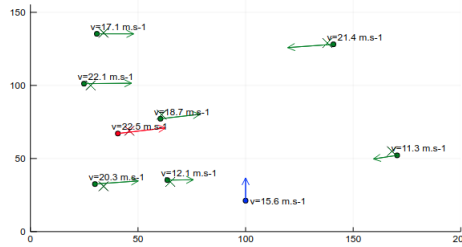
these methods are applicable in restricted cases and usually tested on relatively toy constraints (like a torque constraint). In some cases, when the constraint is not differentiable, a differentiable model is learned to best match the original constraint. The problem we are considering here, in its most generic formulation, is how to best handle complex, non differentiable constraints, in an RL setting. Another key point to consider, is that we have to handle unsafe states, being in an unsafe state (collision risk,  $\min\text{TTC} \leq 10$ ), identifying it as unsafe and trying to move as fast as we can to safer states. In some paper like CPO/TRPO they consider they start in a safe region and try to remain in safe regions while learning new parameters for the Policy Network. So our setting and objectives are different as well: we evaluate a safety cost, if it does not match our hard constraint, we try to improve to match it and then keep it above the hard constraint. So we try to deal with a recovery phase if the hard constraint is not matched. And the main difficulty of the setting is that we are dealing with moving objects for which we have no ground truth models.

We developed an openai gym, module that enable to experiment with some of the challenges.

## 1.2 What did you do?

### 1.2.1 Implementation of an openai gym module

We developed an openai gym test setup for experimentations. We used to have some early code in Julia used in a previous Stanford CS238 course but we migrated everything to Python and properly developed a generic easy to instal and run open ai gym module. It is a test scene where a car, which we call the ego-car, has to drive from point A to point B as fast as possible while othe cars are crossing its path. Other cars may be instantiated randomly at different positions, with different speeds. They may use different driver models: CV, Basic, IDM. It is a model free setting. We are trying to learn an agent that can drive efficiently (as fast as possible from point A to point B with realistic bounded accelerations  $a \in [-2; 2]m s^{-2}$ ) and safely (while minimizing the percentage of collisions).



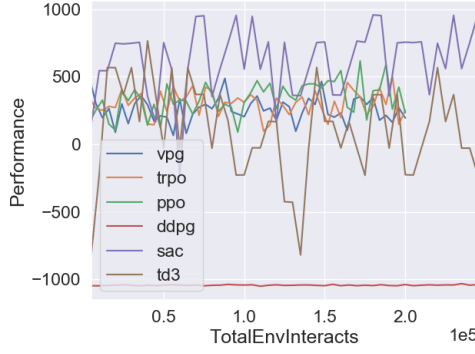
The problem is setup as a MDP with:

- States: different variants will be tested, but as a starting point consider  $\{(x, y, v_x, v_y)_{ego}, (x, y, v_x, v_y)_{obj1..n}\}$
- Actions: longitudinal accelerations  $a_{longi} \in [-2ms^{-2}; +2ms^{-2}]$
- Rewards:  $-1$  for every timestep,  $-1000$  for a collision terminal state,  $+1000$  when goal is reached terminal state
- Model-Free setting: the ego-vehicle does not know the driver model of other cars, it could be CV, Basic or IDM.

An openai gym package has been developed and can be easily used with any existing RL framework. It supports all the standard openai gym API step, reset, render etc and custom ones have been added to experiment with safety constraints and penalties.

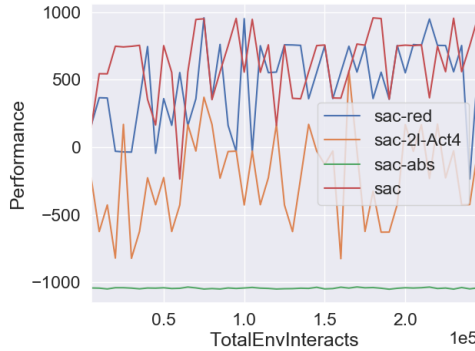
### 1.2.2 Policy Gradients based algorithms benchmark: VPG, PPO, TRPO, DDPG, TD3, SAC

Default setting. Test on Act-v0. SAC performs best.



### 1.2.3 Experiments on different state space representations and how it scale

- With vectorized absolute state space representation: **sac-abs** with a  $\mathbb{R}^9$  state space
  - 1 ego car + 2 cars with  $[x, y, v_x, v_y]$  per car
- With vectorized relative state space representation: **sac** with a  $\mathbb{R}^8$  state space
  - 2 cars relative to ego with  $[x, y, v_x, v_y]$  per car
- With vectorized reduced state space representation: **sac-red** with a  $\mathbb{R}^4$  state space
  - 1 car with minTTC or lowest distance if  $TTC = \infty$
- With vectorized relative state space representation: **sac-2l-Act4** with a  $\mathbb{R}^{16}$  state space
  - 4 cars relative to ego with  $[x, y, v_x, v_y]$  per car
- With image state space representation: could be interesting as independently of the number of cars it will be  $\mathbb{R}^3$  with e.g. 4 stacked images



### 1.2.4 Implementation and Experiments on how to handle complex non differentiable safety constraints and penalties

- Reward shaping with a penalty term
- Line search based on penalty constraint

Backtracking Line Search over a batch of 50000 states tested over different policy networks  $\pi_\theta$

Where  $\theta \leftarrow \theta + \alpha_{bls} \nabla \pi_\theta$  is tested for different step sizes with  $\alpha_{bls} < \alpha_{sgd}$

So we have a 2-stages Gradient Descent, with 2 interleaved steps:

- The first step deals with a differentiable function trying to optimize  $Q(s, a)$  in expectation
- While the 2nd stage is refining the step size  $\alpha$  over a non differentiable safety penalty to search for safer policy networks.

By definition a lot of states are unsafe. They correspond to a predicted  $\min TTC < 10$  seconds. we are trying to come up with a policy that pushes us towards safer states with bounded accelerations at every time step or that ensures we remain in safe states if this is already the case. But we can not move in a single step from an unsafe state to a safer state.

```
Time to compute penalty: 17.25
old_penalty 326022.49
sgd_penalty 322133.54
Average reward: -1881.22 +/- 10.82
Time to compute penalty: 17.11
old_penalty 316778.55
sgd_penalty 354335.09
Backtracking bt_penalty 348790.05
Backtracking: improvement at iter 0 bt_penalty=348790.05 sgd_penalty=354335.09
Backtracking Time: 17.31
Average reward: -1773.23 +/- 9.44
```

**What experiments did you run, what did you find out. It's okay if your new idea wasn't better. But if so, it's important to have ideas about why it didn't work.**

### 1.3 What are next steps / open questions that if you were to continue working, you would do?

- Curse of dimensionality: more efficient state space representation
- Line search is slow: learn where to search first ?

## 2 Introduction

Reinforcement Learning (RL) has demonstrated its capability to learn efficient strategies on many different and complex tasks. In particular, in games like chess and go, the best human players have lost against RL algorithms (Silver et al. [8]). There is a growing traction for applying such RL algorithms to complex robotics tasks like Autonomous Driving. Nevertheless with Autonomous Driving we are dealing with additional challenges. We are in a partially observable environment where enforcing safety is of paramount importance. As a consequence, considering safety via a reward and the optimization of a statistical criteria is not sufficient. Hard Constraints have to be enforced all the time. We propose to study how the RL optimization criteria can be modified to deal with hard constraints; how algorithms like DQN could be modified to cope with such hard constraints and more generally how an RL agent could be integrated in a Decision Making module for Autonomous Driving to provide efficient and scalable strategies while still providing safety guarantees. So we propose to address the following problem formulation:

$$\max_{\theta} \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t, \pi_{\theta}(s_t)) \right]$$

$$\text{s.t. } \text{lower\_bound}(C_i(s_t, a_t)) \geq \text{Margin}_i \forall i \in \llbracket 1, K \rrbracket$$

where the expectation corresponds to the statistical RL objective subject to a set of safety constraints.

We tackle the problem of safe control in physical systems where certain quantities have to be kept constrained. For an autonomous vehicle we must always maintain its distance from obstacles above some margin. But in fact the real state of the world is only partially observable and the Driving Models of surrounding cars are not known exactly: so we are dealing with uncertainty and our constraints are actually a set of random variables  $C_i$  which we want to lower bound. Note that in most of the references the constraints are only considered in expectation via a constraint of type  $J_{C_i}^{\pi} = E_{\pi} [C_i(s, a)]$  whereas here we are interested in enforcing stronger constraints.

**Project Milestone Updates:**

- We will use a single hard constraint which is related to the min Time To Collision that is computed between the ego vehicle and the predicted trajectories of the surrounding vehicles. This min TTC shall be greater or equal than a threshold margin.

### 3 Background/Related Work

In Mirchevska et al. [7] a DQN network is used for tactical decision making in an autonomous driving pipeline but the DQN algorithm itself is not modified to handle hard constraints and the safety is guaranteed by checking the output of the RL algorithm. Our objective here, in contrast, would be to have an RL algorithm that is directly dealing with hard constraints to avoid frequent and sub-optimal actions masking. A review of the different safe RL techniques has been done in García and Fernández [5]. Some techniques mainly deal with soft constraints by either reshaping the reward or trying to minimize the variance related to the risk of making unsafe decisions, while other try to handle hard constraints. Garcia et al. have analyzed and categorized safe RL techniques in two families of approaches: one consists in modifying the exploration process while the other consists in modifying the optimality criterion. In Leurent et al. [6] the RL objective is replaced by a surrogate objective which captures hard constraints and handles model uncertainty by defining a lower bound of the expectation objective. In Achiam et al. [1] constrained policy optimization is solved with a modified trust-region policy gradient. The algorithm’s update rule projects the policy to a safe feasibility set in each iteration. But the policy is kept within constraints only in expectation. In Dalal et al. [4] they directly add to the policy a safety layer that analytically solves an action correction formulation per each state. This safety layer is learned beforehand but is approximated by a first order linear approximator. In Tessler, Mankowitz, and Mannor [9] and in Bohez et al. [2] the proposed approaches are completely in line with our objective here: modifying the RL objective such that it deals directly with hard constraints. But there is no closed form solution for such a problem and a Lagrangian relaxation technique is used for solving the constrained optimization problem. Given a Constrained Markov Decision Process (CMDP), the unconstrained problem is transformed to  $\min_{\lambda \geq 0} \max_{\theta} L(\lambda, \theta) = \min_{\lambda \geq 0} \max_{\theta} [J_R^{\pi_{\theta}} - \lambda(J_C^{\pi_{\theta}} - \alpha)]$  where  $L$  is the Lagrangian and  $\lambda$  the Lagrange multiplier (a penalty coefficient). We propose to study how such techniques could be applied to the Decision Making process of an Autonomous Driving pipeline and we will benchmark different RL algorithms, modified to cope with hard constraints, in an Anti Collision Tests setting.

### 4 Approach

A DQN network will be used as a baseline (we may change later to Policy Gradients or Actor Critic. This is a topic for further refinement). We will consider 3 type of modifications. From the most simple, conceptually, to the most complex, we will:

1. Propose a post DQN safety check. While the DQN network will compute  $Q(s, a_i)$  values for every possible actions, we want to exclude the actions that are unsafe, before deciding what action to take (before taking the  $\operatorname{argmax}_{a_i} Q(s, a_i)$ ). This type of approach is used in a paper from BMW Group by Mirchevska et al. [7].
2. Modify the DQN training algorithm and especially the exploration process so that only safe actions are explored. Similar to Bouton et al. [3], the idea is to derive an exploration strategy that constrains the agent to choose among actions that satisfy safety criteria. Hence the search space of policies is restricted to a “safe” or safer subspace of policies.
3. Replace the RL objective  $\max_{\theta} \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi_{\theta}(s_t))]$  by an objective taking into account hard constraints

$$\max_{\theta} \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi_{\theta}(s_t))]$$

$$\text{s.t. } \text{lower\_bound}(C_i(s_t, a_t)) \geq \text{Margin}_i \forall i \in \llbracket 1, K \rrbracket$$

and study how RL algorithms like DQN should be modified to account for this new objective. In a recent paper from DeepMind from Bohez et al. [2], this type of approach is applied to a realistic, energy-optimized robotic locomotion task, using the Minitaur quadruped developed by Ghost Robotics. But the constraints were considered only in expectation.

### Project Milestone Updates:

- We will focus on step 3
- The surrounding vehicles will move according to an IDM driving model. This driving model depends on 5 parameters. The ego vehicle does not know (precisely) the parameters used by the surrounding vehicles. This is the main source of uncertainty.
- In order to define the hard constraint, in a real use case (not simulated) the ego vehicle shall estimate the Driving Model parameters of every surrounding vehicle with some probability distribution to account for uncertainty. Here in simulation and in the context of this project, we will use some provided range of IDM parameters values: so the IDM model estimation is emulated. Then we will consider worst case scenarios leading to the smallest possible Time To Collision and assess how robust we are against different level of uncertainties.
- So the hard constraint will be  $\min \text{TTC} \geq 10$  where the Time To Collision is computed based on:
  - $\text{acceleration}_{ego-vehicle} = \pi_{\theta}(s_t)$  which depends on the parameters of the Neural Network
  - Predicted trajectories of surrounding vehicles which depend on a set of 5 IDM parameters per surrounding vehicle. We will simulate the fact that these parameters are estimated by the ego vehicle by providing range of values to the ego vehicle in the simulation framework.
  - So the hard constraint will involve the parameters  $\theta$  of  $\pi_{\theta}(s_t)$  such that min Time To Collision is above some threshold.
- In terms of implementing the hard constraint with Tensorflow we will most probably leverage on the following code tensorflow\_constrained\_optimization from Google research.

## 5 Experimental Results

### Project Milestone Updates:

- A github directory has been created for the project: CS234\_Project
- The legacy simulator code has been ported from Julia code (as used in CS238 project) to Python (for use in CS234 Project)
- The simulator code has been upgraded so that the surrounding vehicles move according to an IDM driving model. In previous CS238 project the surrounding vehicles were moving according to a very simple Constant Velocity model.

### 5.1 Simulator

We are upgrading the Anti Collision Tests environment, ACT, developed for a previous CS238 Stanford project. A vehicle has to drive from a point A to a point B as fast as possible, while avoiding other vehicles that are crossing its path and trying to minimize hard braking decisions. So it is a multi-objectives task where efficiency, comfort and safety objectives have to be optimized. While in the previous project we studied the applicability of POMDPs solvers for decision making in a context of sensors uncertainty, we will deal here with an even more challenging task: the uncertainty will be related to the other vehicles driving models. Initially other vehicles driving models were simple Constant Velocity models. Here we will use Intelligent Driver Models, IDM, depending on 5 parameters that will be unknown to the ego vehicle, and that will differ per vehicle. So it is a model-free setting: we do not know the model of the environment, the driving models of others, and we would like to learn to drive efficiently and safely in this context.

### 5.2 Evaluation Metrics

In order to measure success and benchmark different versions of the algorithms, we will use 3 metrics: a safety metric (percentage of collisions), an efficiency metric (time to goal), and a comfort metric (number of hard braking decisions or jerk). We want to enforce safety while not compromising too much efficiency or comfort: a safe AD vehicle that would use many regular hard braking decisions would not be acceptable and could even be dangerous for other vehicles.

## 6 Project Milestones Updates: Remaining Work / Next Steps

- Establish a baseline with RL Policy Gradient algorithm trained on ACT simulator framework. Provide evaluation metrics results.
- Meeting to be planned 5<sup>th</sup> or 6<sup>th</sup> of March with Vaishali, Ramtin, Philippe. Philippe who is located in France will be in California these days.
- Experiment with tensorflow\_constrained\_optimization from Google research on a simple hard constraint use case (e.g. such that the weights of the Neural Network are in a specific range)
- Implement code to compute the hard constraint  $\min \text{TTC} \geq \text{some threshold}$ . In a first step just log how we would deviate from our goal by not enforcing the hard constraint
- Combine Policy Gradient algorithm with the Hard Constraint code during Policy Gradient training
- Compare results with and without enforcing the hard constraint

## References

- [1] Joshua Achiam et al. “Constrained Policy Optimization”. In: (2017). URL: <http://arxiv.org/abs/1705.10528>.
- [2] Steven Bohez et al. *Success at any cost: value constrained model-free continuous control*. 2019. URL: <https://openreview.net/forum?id=rJlJ-2CqtX>.
- [3] Maxime Bouton et al. “Reinforcement learning with probabilistic guarantees for autonomous driving”. In: *Workshop on Safety Risk and Uncertainty in Reinforcement Learning, Conference on Uncertainty in Artificial Intelligence (UAI)*. 2018. URL: <https://drive.google.com/open?id=1d2t14f6GQgH1SERveTmPAR42bMVwHZAZ>.
- [4] Gal Dalal et al. “Safe Exploration in Continuous Action Spaces”. In: (2018). URL: <http://arxiv.org/abs/1801.08757>.
- [5] Javier García and Fernando Fernández. “A Comprehensive Survey on Safe Reinforcement Learning”. In: *Journal of Machine Learning Research* (2015). URL: <http://jmlr.org/papers/v16/garcia15a.html>.
- [6] Edouard Leurent et al. “Approximate Robust Control of Uncertain Dynamical Systems”. In: 2018.
- [7] Branka Mirchevska et al. “High-level Decision Making for Safe and Reasonable Autonomous Lane Changing using Reinforcement Learning”. In: 2018. DOI: 10.1109/ITSC.2018.8569448.
- [8] David Silver et al. “Mastering the game of Go without human knowledge”. In: *Nature* (2017). URL: <http://dx.doi.org/10.1038/nature24270>.
- [9] Chen Tessler, Daniel J. Mankowitz, and Shie Mannor. “Reward Constrained Policy Optimization”. In: abs/1805.11074 (2018). URL: <http://arxiv.org/abs/1805.11074>.