Improvement of Vehicle Stability Using Reinforcement Learning

Janaína R. Amaral^{1, 2}, Harald Göllinger², Thiago A. Fiorentin¹

¹ Pós-Graduação em Engenharia e Ciências Mecânicas – Universidade Federal de Santa Catarina (UFSC) Joinville – SC – Brazil

²Department of Mechanical Engineering – Technische Hochschule Ingolstadt Ingolstadt, Germany

janainaribasdeamaral@gmail.com, Harald.Goellinger@thi.de, thiago.antonio@ufsc

Abstract. This paper presents a preliminary study on the use of reinforcement learning to control the torque vectoring of a small rear wheel driven electric race car in order to improve vehicle handling and vehicle stability. The reinforcement learning algorithm used is Neural Fitted Q Iteration and the sampling of experiences is based on simulations of the vehicle behavior using the software CarMaker. The cost function is based on the position of the states on the phase-plane of sideslip angle and sideslip angular velocity. The resulting controller is able to improve the vehicle handling and stability with a significant reduction in vehicle sideslip angle.

1. Introduction

Controllers based on reinforcement learning algorithms, which are widely used to control the motion of mobile robots [Benbrahim and Franklin 1997] [Duan et al. 2008] [Janusz and Riedmiller 1995] [Riedmiller et al. 2009] [Vincent and Sun 2012], are gaining space in the automotive area. It was already applied for road following by controlling the steering angle [Oh, Lee and Choi 2000] [Riedmiller, Montemerlo, Dahlkamp 2007] [Yu and Sethi 1995], for enhancing ride performance by controlling the suspension system [Howell et al. 1997] and for improving longitudinal motion in Adaptive Cruise Control [Desjardins and Chaib-draa 2011] [Pietquin, Tango and Aras 2011]. It was also used to perform lane change autonomously [Wang, Chan and de La Fortelle 2018] and regarding lateral vehicle dynamics, a reinforcement learning algorithm was used to optimize rules of a fuzzy system that controls vehicle stability adding a yaw moment generated by differential braking [Akbari and Goharimanesh 2014].

In this work, another application of reinforcement learning (RL) in the automotive area is studied. Here, a reinforcement learning algorithm called Neural Fitted Q Iteration [Riedmiller 2005] is used to control the torque vectoring of a small rear wheel driven electric race car in order to improve the vehicle handling and stability.

As reinforcement learning algorithms learn how to take decisions by trial and error, the algorithm learned how to control the torque vectoring experiencing interactions with the environment by means of simulations of the vehicle behavior that are carried out using the software Carmaker. The manoeuvre selected was Sine with Dwell which is used by UN ECE for the approval of ESC in passenger cars. However, it was necessary to set costs (penalties) for unstable or handling limit situations to help the algorithm understand when it was necessary to add a yaw moment to keep the car stable.

One of the most important sets for the stability controllers is to define the threshold for the addition of the corrective yaw moment. Some controllers compare the current yaw rate of the vehicle with the desired yaw rate calculated using a linear or nonlinear vehicle model [Goharimanesh and Akbari 2017] [Lee, Hwang and Suh 2015]. Another approach is to define the handling limit using the phase-plane of sideslip angle and sideslip angular velocity. By looking at this phase-plane it is possible to have a good characterization of the nonlinear dynamic behavior of the vehicle. For this reason, it is used in many stability controllers that use brakes, torque vectoring or even a combination of both solutions [Guo et al. 2010] [Inagaki, Kshiro and Yamamoto 1994] [He 2005] [Lu et al. 2016]. In this work, the phase-plane of sideslip angle and sideslip angular velocity is used to define the cost function.

At the end, a preliminary handling and stability controller is created and it is able to select which percentage of torque should be distributed to the wheels when a steering input as sine with dwell is given.

Therefore, in this article a brief introduction about lateral vehicle dynamics and reinforcement learning is made. In the sequence, the process to create the controller is presented. Then, the learning results are shown and discussed.

2. Lateral Vehicle Dynamics and Handling Limit

When the driver needs to turn the vehicle in a situation of cornering or lane change, for example, he applies a steering angle input to the vehicle. According to this input, the wheels are turned into the desired direction.

If the vehicle is at a high speed, the tires must develop lateral forces to counteract the lateral acceleration that is present under this situation to control the direction of the vehicle. The generation of the lateral force is based on lateral slip of the tire (slip angle), on lateral inclination (camber angle) or the two effects combined [Gillespie 1992].

The slip angle α is the angle between the direction of heading of the tire and its direction of travel [Gillespie 1992], as can be seen in Figure 1. The lateral force generated by the tire is dependent on the slip angle. For small values of sideslip, the relation is linear. However, when the slip angle increases and the lateral acceleration exceeds 0.3g, the relation becomes nonlinear [He 2005], as shown in Figure 1.

After further increase in the side slip motion, the properties of the lateral tire forces saturate [He 2005] and the tire presents the behavior of a locked wheel [Gillespie 1992]. The result is a yaw moment that leads to vehicle instability and spin [He 2005]. From this point, the vehicle will not respond to steering inputs, therefore it is the limit of handling, once that handling is the responsiveness of a vehicle to driver input [Gillespie 1992].



Figure 1. Lateral force and slip angle. Source: [Gillespie 1992].

To control the vehicle and bring it back to a stable state, a yaw moment can be added by means of braking one of the wheels or by applying torque vectoring (setting different values of torque for the wheels). Using the brakes, it is possible to have a higher yaw moment, but it affects the longitudinal motion of the vehicle. In the case of controllers that apply torque vectoring, they do not affect the longitudinal dynamics and they are effective throughout the handling regime [He 2005].

3. Reinforcement Learning (RL)

Reinforcement learning is a machine learning approach to solve problems in which an agent needs to learn the optimal strategy of actions interacting with the environment by means of experiences of trial and error. During the process of learning, the agent observes states, selects actions and receives rewards or costs for the transition state-action [Sutton and Barto 1998] [Wiering and Otterlo 2012].

The rewards represent the desirability of the state in an immediate sense. The goal of the RL algorithms is to maximize the sum of rewards or minimize the sum of costs [Sutton and Barto 1998]. To achieve the goal, the agent needs to learn the optimal policy, that is the optimal way of selecting the actions.

If the decision problem can be considered a Markov Decision Problem, it is described by a set of states s, a set of actions a, the probability of the transitions to a new state s' p(s, a, s') and the immediate reward or cost function [Sutton and Barto 1998]. Then it can be solved by looking at the state value function or the action value function. The state value function represents the expected amount of rewards that can be reached from one state, in other words, it tells how good is to be in one state [Wiering and Otterlo 2012]. The action value function (Q-value) represents the expected amount of reward that can be reached by being in one state and take one action.

When the probability of the transitions is known, these values are easily calculated by dynamic programming: looking at all states and using the Bellman Equations iteratively. But, when the probability is unknown it is called model-free approach and the value functions are calculated by interacting with the environment and collecting experience samples [Sutton and Barto 1998].

At the end, the optimal policy can be found by selecting the action that leads to the next state with highest state value or to a next transition state-action with highest action value. However, when the state space is large and complex, a lot of memory would be necessary to store all the information. In this case, the Q-value can be approximated by a function that normally is represented by a neural network or a decision tree [Sutton and Barto 1998].

In the model-free approach, many methods can be used to find the optimal policy, nevertheless when function approximation is necessary, the batch reinforcement learning is more appropriate due to its offline update of the Q-value [Lange, Gabel and Riedmiller 2012]. The Neural Fitted Q Iteration (NFQ) is an example of batch algorithm that uses a multi-layer perceptron as Q-value function [Riedmiller 2005].

3.1. Neural Fitted Q Iteration

The NFQ algorithm is efficient and effective regarding the training of the Q-value function as it is made after a set of experiences are stored [Riedmiller 2005]. The idea of the algorithm is to generate the training data set P from a set of transitions experiences called sample set D collected in the form (s, a, s') and then, to train of the neural network to update the value function.

The training data set P must be composed by the input that is the pairs state-action seen in the sample set D and the target that is calculated by adding the transition cost with the lowest Q-value for the next state s' obtained using the current Q-value function [Riedmiller 2005]. Once that P is generated, the training of a new neural network is performed and the new replaces the old one.

Later, when the process is finished, more data can be collected following a policy greedy (using the neural network to select the action related to the lowest Q-value) or ε -greedy (alternating between greedy selection and random selection). In the sequence, the new data is added to the sample set D and more NFQ iterations can be carried out to update the Q-value function.

4. RL Controller Design

The goal of the proposed controller is to observe the current state of the vehicle, which is provided by sensors, and define the percentage of torque that should be delivered to each wheel to generate an additional yaw moment in cases when the vehicle is in handling limit or unstable situation.

Different from traditional controllers for stability that depends on a mathematical model of the vehicle, the RL controller learnt the torque distribution by interacting with the environment observing states, taking actions and receiving penalties.

4.1. Controller Inputs

In order to select the input states, simulations with the race car were performed in the software CarMaker. During the simulations, states that represent the longitudinal and lateral behavior of the vehicle and input of the driver were collected: sideslip angle in rad, sideslip angular velocity in rad/s, speed in m/s, yaw rate in rad/s, steering angle in rad, acceleration in x and in y in m/s² (longitudinal and transverse respectively).

With the collected data it was possible to see that steering angle, yaw rate and acceleration in y are highly correlated, but the correlation among inputs is not recommended when neural networks are used because it might slow the learning [LeCun

et al. 1998]. Then, the steering angle and the yaw rate where selected as inputs because they are used to describe the lateral vehicle dynamics [Gillespie 1992] and these two variables were decorrelated using Principal Component Analysis (PCA). It resulted in two new variables PC1 and PC2 that are a linear combination of the yaw rate and steering angle.

The speed and the acceleration in x (in longitudinal direction) were also selected because these variables are easily measured. On the other hand, sensors for sideslip angle and sideslip angular velocity normally are expensive, so these variables were not considered as an input in this work. It makes the application of the algorithm simpler. Therefore, the selected inputs were: speed, PC1, PC2 and acceleration in x. The value of the inputs was standardized.

4.2. Possible Actions

The actions in each state were defined as the percentage of torque that goes to the left rear wheel and the possible percentages were 30%, 40%, 50%, 60% and 70%.

As the vehicle has two driven wheels but it has only one motor, it is necessary that the differential distribute different torque to the wheels.

4.3. Cost Function

The cost function was based on the analysis of the phase-plane of the vehicle sideslip angle β (between the speed and the longitudinal axis) and sideslip angular velocity β ' proposed by [He 2005]. The phase-plane is a graphical method used to perform the stability analysis of nonlinear dynamic systems and it is obtained by plotting one state as function of the other. At the end it represents the transient response of the system to initial conditions or constant inputs [He 2005].

In the case of the vehicle, its lateral vehicle dynamics can be represented by a second order and nonlinear system [Gillespie 1992]. Therefore, it is possible to use the phase-plane diagram to perform the stability analysis of vehicles. However, to create the phase-plane diagram of the vehicle, it is important to know the vehicle model and the vehicle tire model. But as this diagram was used only to evaluated how good or how bad the actions of the RL algorithm were, in this work the phase-plane for the studied vehicle was not created. Instead, the β - β ' phase-plane from [He 2005] was used.

[He 2005] proposed the division of the phase-plane diagram in 3 regions: reference region (stable situation), stability error (stable situation or handling limit) and unstable region as presented in Figure 2.



Figure 2. Definition of reference region and stability error for vehicle control. Source: [He 2005].

In the present work, this idea was used to create the vehicle cost function. It was considered that when the vehicle was in state s, took action a and achieved a state s' that is in region 3 (emergency situation), it was a failure. In this case, the value of the cost function is 1. If the state s' is in the "stability error" region, the value was 0.4. When the state s' is in the reference region (stable situation) and the action taken is 50% of torque distribution, the value of the cost function for the state-action pair is zero. But when the vehicle is in a stable situation and the torque distribution was different than 50%, the cost value is 0.01. This consideration was done because the vehicle should learn not to use torque distribution different than 50% in stable situation as it can make the handling unnatural for the driver [He 2005]. The cost function using the regions defined by [He 2005] can be seen in Equation 1.

$$c(s, a, s') = \begin{cases} 1.00, & if |\dot{\beta} + 4\beta| \ge 72 \\ 0.40, & if 24 \le |\dot{\beta} + 4\beta| < 72 \\ 0.01, & if |\dot{\beta} + 4\beta| < 24 \& action \ne 50\% \\ 0.00, & if |\dot{\beta} + 4\beta| < 24 \& action = 50\% \end{cases}$$
(1)

4.4. Q-value Function

In the Neural Fitted Q Iteration, a regression neural network is used to estimate the Qvalue, which represents how good is to be in one state and take some action. The neural network architecture used was composed of one input layer with five nodes, two hidden layers with ten nodes and output layer with one node. For the nodes of the hidden layers the activation function was the logistic sigmoidal (0,1) as used for [Riedmiller, Montemerlo, Dahlkamp 2007], and for the output layer was the function purelin. The Scaled Conjugate Gradient algorithm was used for the supervised training of the neural network with a maximum of 400 epochs. The input to train the neural network was the state observed (speed, PC1, PC2, and longitudinal acceleration) and the action taken in that state. The target was the estimated Q-value for the pair state-action. It was calculated according to Equation 2, where the cost of the transition is added to the minimum Q-value related to the next state (calculated using the previous neural network) multiplied by the discount factor γ considered as 0.95. Each time a new neural network was trained, the estimated Q-value was improved.

$$Q(s,a) = c(s,a,s') + \gamma \min_{a'} Q(s',a')$$
(2)

4.5. Policy

As the data collection to train the NFQ algorithm was performed by means of simulations in CarMaker, it was possible to use the policy ε -greedy with 10% of random choices to explore the state space and 90% of greedy choices (selecting the action with lowest Q value in the state) to improve the policy.

4.6. Sampling of Experiences

The sampling of experiences was performed by means of simulation of the dynamic behavior of a small rear wheel driven race car (similar to a Formula Student race car) in the software Carmaker. This race car has a weight of 191 kg, 1.6 m of wheelbase and track width equals to 1.2 m.

Running the CarMaker in Simulink, it was possible to add a controller for torque vectoring in the vehicle model in software-in-the-loop configuration. Therefore, the controller could read the states (controller inputs), calculate the Q-values using one neural network block for each state-action pair and take actions according to the policy ε -greedy.

The situation analysed was the manoeuvre Sine with Dwell used by UN ECE R13H [UN ECE 2014] for the approval of the Electronic Stability Control in passenger cars. In Figure 3, it is possible to see the shape of the steering wheel angle used as input.



Figure 3. Steering input used in the manoeuvre Sine with Dwell. Source: [UN ECE 2014].

According to UN ECE R13H, to approve the ESC it is necessary to find firstly the steering angle A that causes in the vehicle a lateral acceleration of 0.3g at 80 km/h. In the

sequence, it is necessary to perform simulations at 80 km/h with the amplitude of the sine input being the value A multiplied by factors ranging from 1.5 to 6.5 in steps of 0.5 (clockwise and anticlockwise) and then, verify if they comply with the criteria defined by the standard. In this work, due to limitation of computational time, the simulations for sampling of experiences were performed for 5.5A (63.23°) because it is the smallest amplitude that does not comply with the criteria; 6.5A (74.73°) because this was the highest factor that should be used in the simulation; 8A (91.97°) because, according to the phase-plane of sideslip angle and its angular velocity, this is the smallest angle that leads to loss of stability.

As the vehicle used in the simulations was smaller than a passenger car, the amplitudes for the sine input were also smaller than what would be expected for a passenger car. However, this standard was still used because of its complete procedure and adaptation of the amplitudes regarding the characteristics of the vehicle.

4.7. Learning Process

The learning process used in this work is represented by the alternation of steps of sampling of experience with steps of updating of the Q-value function, because the first neural network must be trained with random values and does not represent the Q-value. Therefore, to improve the estimation of the Q-value, after each sampling performed in CarMaker, the sample data set obtained was added to a memory and one iteration of NFQ was carried out in Matlab to generate the training data set P and to train a new neural network. The new neural network was replaced in Simulink and another step of sampling of experience was performed. This cycle is shown in Figure 4.



Figure 4. Representation of the learning process.

This process was repeated eight times for the six variations of sine amplitudes $(\pm 63.23^\circ, \pm 74.73^\circ, \pm 91.97^\circ)$ and for the first sampling of experiences a neural network trained with random values in a range of [0, 1.5] was added to the vehicle model in Simulink.

It is important to point out that to have a good controller it is necessary more simulations for sampling of experiences, this is only the preliminary study.

4.8. Final Controller Architecture

After the algorithm learnt successfully the behavior of the system, it was able to control vehicle handling and vehicle stability by receiving information of the state, using neural networks to calculate the Q-values for all state-action pairs. Then, the action with the lowest Q-value is selected. In other words, the percentage of torque that should be distributed to each wheel to keep the car stable. For that, five identical neural networks were used, one for each possible action. The input of the neural networks is the state s and the defined action. This process is represented in Figure 5.



Figure 5. Representation of the final controller.

5. Results of the Learning Process

At the end of the learning process, the performance of the preliminary controller was evaluated by means of comparison between the result of simulations of the vehicle behavior with and without controller.

The simulations were carried out in CarMaker and the manoeuvre performed was Sine with Dwell with amplitude 63.23 degrees and 91.97 degrees at 80 km/h. In Figure 6, the results of the simulations are presented: the steering input, the yaw rate and the phase-plane of the vehicle sideslip for both amplitudes.

In Figure 6.(e), it is possible to see that without the controller the vehicle would not leave the regions one and two of the phase-plane. But with the controller, the peak of the vehicle sideslip angle was reduced in 34,60%. Therefore, vehicle handling was improved.

Figure 6.(f) shows that without the controller the vehicle loses stability as the β - β ' phase trajectory reaches region 3. With the controller, the vehicle remains within the limits of handling, which means stability enhancement. It can also be seen in Figure 6.(c) and 6.(d), because after that the steering wheel input ends at 41.9 s, the yaw rate stabilizes faster using the controller.

Moreover, comparing the steering input with the yaw rate for each simulation, it is possible to say that the controller was able to improve the tracking.



42.5

43



10

15

0 -50

-100

-150

-200

-20

-15

-10 -5 0

5

10 15 25

30

20

6. Conclusions

-10

-5

0

5

0

-50

100

-150 -15

This paper presents the use of reinforcement learning to control the torque vectoring of an electric vehicle in order to improve vehicle handling and vehicle stability. Preliminary results of the learning process show that the selected cost function and states were satisfactory for learning to cope with the vehicle behavior. The results also show that the RL controller could improve the vehicle dynamics and the vehicle's ability to follow the track.

To improve the controller performance, more simulations for sampling of experiences are going to be done. Also, simulations with other steering inputs and other manoeuvres could be added to the learning process. Furthermore, it could be considered different percentages of torque distribution to allow the addition of higher yaw moments.

7. Acknowledgment

This paper was a result from the cooperation between THI and UFSC by means of the German-Brazilian Academic Project Applied Network on Automotive Research and Education (AWARE). The authors thank Fundação de Amparo à Pesquisa e Inovação do Estado de Santa Catarina (FAPESC) for financial support.

References

- Akbari, A. A. and Goharimanesh, M. (2014). Yaw Moment Control Using Fuzzy Reinforcement Learning. In *International Symposium on Advanced Vehicle Control*.
- Benbrahim, H. and Franklin, J. A. (1997). Biped dynamic walking using reinforcement learning. In *Robotics and Autonomous Systems 22*, pages 283-302. Elsevier.
- Desjardins, C. and Chaib-Draa, B. (2011). Cooperative adaptive cruise control: A reinforcement learning approach. In *IEEE Transactions on intelligent transportation systems*, pages 1248-1260. IEEE.
- Duan Y., Cui B. and Yang H. (2008). Robot Navigation Based on Fuzzy RL Algorithm. In: Sun F., Zhang J., Tan Y., Cao J. and Yu W. (eds) *Advances in Neural Networks*. Lecture Notes in Computer Science, vol 5263. Springer, Berlin, Heidelberg.
- Gillespie, T. D. (1992), Fundamentals of Vehicle Dynamics, Society of Automotive Engineers.
- Goharimanesh, M. and Akbari, A. A. (2017). Improving lateral dynamic of vehicle using direct yaw moment controller by differential brake torques based on quantitative feedback theory. In *Scientia Iranica. Transaction B, Mechanical Engineering*, pages 662-672.
- Guo, J., Chu, L., Liu, H., Shang, M. and Fang, Y. (2010). Integrated control of active front steering and electronic stability program. In *International Conference on Advanced Computer Control* (ICACC), pages 449-453. IEEE.
- He, J. (2005). Integrated vehicle dynamics control using active steering, driveline and braking (Doctoral dissertation, University of Leeds).
- Howell, M. N., Frost, G. P., Gordon, T. J. and Wu, Q. H. (1997). Continuous action reinforcement learning applied to vehicle suspension control. In *Mechatronics*, pages 263–276. Elsevier.
- Inagaki, S., Kshiro, I. and Yamamoto, M. (1994). Analysis on vehicle stability in critical cornering using phase-plane method. In *International Symposium on Advanced Vehicle Control*, pages 287-292.
- Janusz, B. and Riedmiller, M. (1995). Self-learning neural control of a mobile robot. In *IEEE International Conference on Neural Networks*, Vol. 5, pages 2358-2363. IEEE.

- Lange, S., Gabel, T. and Riedmiller, M. (2012) "Batch Reinforcement Learning", In: M. Reinforcement Learning. Adaptation, Learning, and Optimization, Edited by Marco Wiering and Matijn van Otterlo, vol 12, Springer, Berlin, Heidelberg.
- LeCun, Y., Bottou, L., Orr, G. B. and Muller, K.R. (1998). Efficient BackProp. In *Neural Networks: Tricks of the trade*. Springer.
- Lee, M., Hwang, K. and Suh, I. S. (2015). Independent and Integrated Torque Control of 4-Wheel Drive Electric Vehicle for Automated Driving. In *International Electric Vehicle Symposium and Exhibition*.
- Lu, Q., Gentile, P., Tota, A., Sorniotti, A., Gruber, P., Costamagna, F. and De Smet, J. (2016). Enhancing vehicle cornering limit through sideslip and yaw rate control. In *Mechanical Systems and Signal Processing*, pages 455-472. Elsevier.
- Oh, S. Y., Lee, J. H. and Choi, D. H. (2000). A new reinforcement learning vehicle control architecture for vision-based road following. In *IEEE Transactions on Vehicular Technology*, pages 997-1005. IEEE.
- Pietquin, O., Tango, F. and Aras, R. (2011). Batch reinforcement learning for optimizing longitudinal driving assistance strategies. In *IEEE Symposium on Computational intelligence in vehicles and transportation systems*, pages 73-79. IEEE.
- Riedmiller, M. (2005). Neural fitted Q iteration-first experiences with a data efficient neural reinforcement learning method. In: *European Conference on Machine Learning*, pages 317-328. Springer, Berlin, Heidelberg.
- Riedmiller, M., Montemerlo, M. and Dahlkamp, H. (2007). Learning to drive a real car in 20 minutes. In *Frontiers in the Convergence of Bioscience and Information Technologies*, pages 645-650. IEEE.
- Riedmiller, M., Gabel, T., Hafner, R., and Lange, S. (2009). Reinforcement learning for robot soccer. In *Autonomous Robots*, pages 55–73. Springer.
- Sutton, R. S. and Barto, A. G. (1998), Reinforcement Learning: an Introduction, MIT Press.
- UN ECE (2014). Addendum 12-H: Regulation No. 13-H, third edition.
- Vincent, I. and Sun, Q. (2012). A combined reactive and reinforcement learning controller for an autonomous tracked vehicle. In *Robotics and Autonomous Systems*, pages 599-608. Elsevier.
- Wang, P., Chan, C. Y. and de La Fortelle, A. (2018). A Reinforcement Learning Based Approach for Automated Lane Change Maneuvers. In 2018 IEEE Intelligent Vehicles Symposium. IEEE.
- Wiering, M. and Van Otterlo, M. (2012), Reinforcement learning. Adaptation, learning, and optimization, Springer, v. 12.
- Yu, G. and Sethi, I. K. (1995). Road-following with continuous learning. In *Intelligent Vehicles' 95 Symposium*, pages 412-417. IEEE.