

Solving the Inverted Pendulum Control Problem using Reinforcement Learning

Aalap Doshi

*Electrical and Computer Engineering
Arizona State University*

Poojan Patel

*Electrical and Computer Engineering
Arizona State University*

Varadaraya Ganesh Shenoy

*Electrical and Computer Engineering
Arizona State University*

Abstract—The inverted pendulum is a typical, intriguing control issue that includes numerous essential components of control hypothesis. It is a classical benchmark control problem because its dynamics resembles with that of numerous real world systems of interest like missile launchers, human walking, segways and many more. The control of this system is challenging as it is exceptionally unsteady, profoundly non-straight, non-minimum phase system and under actuated. Also, there are physical constraints on the track position control voltage complexity in its control design. We need a controller to balance the pendulum such that the cart on the track can be controlled quickly and accurately so that the cart remains in the center and the pendulum is always straight in its inverted position during its motion. Reinforcement learning methods are discussed which follows with the idea of punishing for bad actions and rewarding for good ones.

Index Terms—Inverted Pendulum, Reinforcement Learning, Temporal Difference

I. PROBLEM STATEMENT

The inverted pendulum is a popular research problem as it is an unstable system (ie) the pendulum will fall over if the movement of the cart doesn't balance it. Thus, the primary objective of the control design is to apply a force to the cart in order to balance the pendulum.

However, orthodox design techniques are successful given knowledge of the system including its dynamics, and an objective function expresses the desired behavior of the system. But, how can control be achieved if this knowledge is not available? This question is solved by learning a function which selects control actions given the current state of the pendulum, through experience of trying various actions and examining the results, with no history pertaining to the correct action to begin with.

This type of controller design where the objective function which evaluate states and actions is absent, can be modified only on basis of occurrence of failure signals. This results in an assignment-of-credit problem, which is necessary to trace the sequence of actions leading to failure.[And89]

In this project, a 2D problem in which the pendulum movement is restricted along the vertical plane. The control input for the system is the force responsible for horizontal movement of the cart while outputs are the angular position and the cart position.

The project implements machine learning methods that help realize successful action sequences by generating two functions, namely action and evaluation functions respectively. Current state of the system is translated to the corresponding control using action function while the mapping of the current state into its evaluation is carried out by the latter. These functions are learnt using two networks referred to as : Action Network and Evaluation Network.

II. INVERTED PENDULUM

In this case we will consider a two-dimensional issue where the pendulum is compelled to move in the vertical plane appeared as depicted in the figure. For this framework, the control input is the force F_t that moves the cart horizontally on the level plane and the

outputs are the angular position of the pendulum θ_t and the horizontal position of the cart h_t .

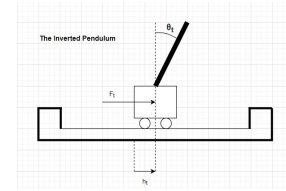


Fig. 1: The Inverted Pendulum

The θ_t and its state is used to balance the pendulum while the cart is moving in the horizontal direction both in right and left direction.

Initially we observe a lot of change in the angle because of the disturbance and constant movement of cart. As time moves forward we see that the error in the angle reduces as the systems starts learning from the failure signal.

Now, what is failure signal? Here, failure signal can be described in two parts: first being the change in the angles in degree and second the horizontal movement of cart in both the direction. We can always assign a threshold for the system for both θ_t and h_t so that when the system starts balancing the inverted pendulum with each motion of pendulum and cart the system learns from its failure to balance it perfectly with each attempt and manages to reduce the error every time.

For example, the pendulum falling past +15 degree or the cart hitting the bounds of the track at +3.0m. The goal of the inverted pendulum task is to apply a sequence of right and left forces of fixed magnitude to the cart such that the pendulum is balanced and the cart does not hit the edge of the track.

The objective as simply expressed makes this errand exceptionally troublesome; the failure signal is a delayed and rare performance measure. Before portraying an answer for this plan of the altered pendulum task, we quickly examine different methodologies that assume the existence of additional task-specific knowledge.

Here, we try the approach with state equation and state variables.

III. FORCE EQUATIONS AND SYSTEM ANALYSIS

State variable of the system are those parameters of interest which can be utilized to locate every single other parameter of the system and whose information permit us to think about the present or future condition of the general system. And State Equation shows the relationship between the system's actual state and its input, and the future state of the system. The Output Equation shows the relationship between the system state and its input, and the output.

We tested the approach which used just the state variables and state equation along with some motion equations to control motion of inverted pendulum.

According to the free-body diagrams of the two elements of the system as illustrated below [Mic],[Cannon] :

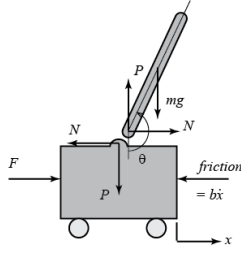


Fig. 2: Free-Body Diagram

Summing the forces in horizontal and vertical directions of the cart, the equation obtained for the horizontal direction is,

$$M\ddot{x} + b\dot{x} + N = F \quad (1)$$

The following equation is got when the forces in horizontal direction of the pendulum is summed,

$$N = m\ddot{x} + ml\ddot{\theta}\cos(\theta) - ml\dot{\theta}^2\sin(\theta) \quad (2)$$

Putting (2) in (1),

$$(M + m)\ddot{x} + b\dot{x} + ml\ddot{\theta}\cos(\theta) - ml\dot{\theta}^2\sin(\theta) = F \quad (3)$$

When the forces perpendicular to the pendulum are added , the second equation of motion is derived,

$$P\sin(\theta) + N\cos(\theta) - mg\sin(\theta) = ml\ddot{\theta} + m\ddot{x}\cos(\theta) \quad (4)$$

In order to eliminate P and N in the above equation, the centroids about the pendulum is added,

$$-Pl\sin(\theta) - Nl\cos(\theta) = I\ddot{\theta} \quad (5)$$

Combining (4) and (5) we get,

$$(I + ml^2)\ddot{\theta} + mgl\sin(\theta) = -ml\ddot{x}\cos(\theta) \quad (6)$$

It is observed that eq(6) represents a non-linear equation, which calls for linearization of the system. For the inverted pendulum system, it is realistic to linearize the system along the vertical axis. Thereby, $\theta = \pi$ is the equilibrium chosen assuming the system will always be in the neighborhood of this criteria. Thus, resulting in θ becoming $\pi + \phi$.

The traditional approach for the control and modeling of any system is when we know all the dynamics of the system and all the motion equations of the system.

For inverted pendulum problem the equations are assumed that the control Force F is a linear function of the four state variables $(\theta_t, \dot{\theta}_t, \theta_t, \dot{\theta}_t)$ and in an equation form it would have constant coefficients multiplied with it. From the exercise carried out and also when we read about how state equations and state variables can be used for dynamic control of inverted pendulum and also cart motion. It was found that this approach was tested by [CL87] in which they used linear feedback in three out of four variables to test stable control of a ball-balancing experiment. The experiment's success heavily depends on state equations and the linearized approximation.

However, we need a compact system which is able to control in any dynamical situation and it does not have any limitation regarding any environmental conditions for balance. We found that the first use of neural network was performed using state variables and equations and then after seeing the input-output behavior of the linear control

systems mimic the same system by training the neural network. This caused many different approaches to arise with different limitations.

As per our analysis multilayer networks can also be used with nonlinear control laws but here as well the system dynamics should be known. What if we do not know the system dynamics but we still want the desired result? Here, we would need some adaptive or learning approach to obtain a stable control. Even a novel approach of completely eliminating state variables and using image processing to derive visual image of the pendulum location and giving a real time feedback to the system can be used. But, all this approach cannot give us an optimal solution.

One of the interesting solutions we found was that if neither a designed controller nor a human expert is available we should make learning based on the actual parameters and keep the feedback loop closed so that each and every time an error is generated it learns from that and reduces it in the next go. A reference states can be assigned for example $(0, 0, 0, 0)$ and based on difference from the current state to reference state each action can be taken to reduce this difference. All the above discussion is effective but has some or other limitation which eliminates its use.

Following we discuss neural network approach using networks, state variables and machine learning technique to solve the problem.

IV. DISCUSSION AND LITERATURE REVIEW

A. Single Layer Approach and Reinforcement Learning

In this section we will discuss about a different approach in which we learn about neural network approach and learning methods such as supervised learning, unsupervised learning. Supervised learning methods are most commonly used learning method in neural network but it requires a training data consisting of input data and output data.

For the inverted pendulum problem we are not aware about the output data since it varies with environment so we cannot predict the correct action. Here, Reinforcement learning comes into the picture because it learns from the effects measured by changes in an evaluation signal. While we do not know the output vectors we can complete the loop by checking each failure signal and reduce the error by each iteration. The method which we used has two single layer network with one named action network and other being evaluation network.

[And89] Action network penalizes each control action of the pendulum for maintaining the angle and for cart to push left and right according the magnitude of force. Output of this unit is probabilistic which means the initial weights are zero and then the action units learn via a reinforcement learning method. With each movement it learns and assigns different weights to the system. So, even with a slight movement learning action takes place but it is very slow because of the delayed signal.

Second network is called the evaluation network, which also consists of single unit. This unit specially calculates the Temporal Difference (TD) which measures the sum of future failure signals by prediction. The prediction is based on the output of the inverted pendulum state vectors and failure signal. Through learning, the output of the evaluation network comes to predict the failure signal, with the quality of the expectation showing how soon disappointment can be relied upon to happen. With each step of prediction of network's input and the change in the value of new prediction (the difference), based on the current state of the inverted pendulum, and the previous prediction, based on the previous state the predictions are adjusted accordingly. This sequence of event goes on and changes ends with occurrence of failure signal.

The final prediction is based on failure signal and the previous prediction. This is the reason why evaluation network is very important. Inner product of the input vector and unit's weight vector gives output. Input vectors are the states of inverted pendulum and weights are developed by the TD method. With each ranking to the states the difference in the units output on the transition from one state to another is used to judge effectiveness of previous action. A decrease in failure signal proves that the evaluation is effective and probability of prediction is increased.

In this way reinforcement learning is performed. However, one problem which we found is that learning is not sufficient in initial stage without much experience so the weights updated at that time will be very uncertain and that causes disturbance in the system initially.

In this particular case, four state variables are adequate to represent action unit since the output for it is linear. The failure signal is given value -1 and 0 to other states. Failure occurs when the value of θ above or below the deviation ϕ . Evaluation unit is basically function of just θ (pendulum's angle) and $\dot{\theta}$ (angular velocity). With the pendulum moving towards balanced position, evaluation is closer to 0 which means a weaker prediction of failure. But we need to develop a different representation of function for evaluation unit.

B. Advantages and Disadvantages of the Single Layer Approach

We found a method an experiment which was motivated by the above-discussed methods. [MC68] in there experiment devised a system called "BOXES" which learned to control an inverted pendulum using discretized state-space representation. The four state variables $\theta_t, \dot{\theta}_t, h_t, \dot{h}_t$ are intersected and form 162 regions.

As discussed above, there are two networks (action and evaluation) and each receives 162 binary input components. All the process discussed above for each network is performed but the key difference they used which we can also implement is that they used counters in each region to store and remember the past states and failure. Hence, learning occurs only on a failure signal.

The TD method allows learning based on the continuous learned evaluation function and uses the output of that network as reinforcement as difference. With this method, we found that the probability that the network will generate actions leading to failure tends to zero and the probability of learning and balancing increase with each iteration. Advantage of this method is that it does not depend upon the future failures to learn but it learns from the difference in evaluation function and its output. This is an advantage because the probability of failure reduces and it does not depend upon the summation of failure signals but depends on each output generated and looped back.

One disadvantage which we came through reading many resources online is that to divide the state space, one must strike a balance between generalizing and learning speed. Quantization can be performed for all complex function but performing quantization for all regions requires much experience. Solution for this is that learning can be made faster using coarse quantization in which learning from one state is transferred to all states in the region. And, here the function whose output remains relatively constant over regions can be represented. This is an adaptive representation that learns the features based on the experience.

C. Solution using two-layer networks

Two-single layer network has some limitations like assigning state variables with zero weight on every failure will balance the pendulum but will not keep the cart at the center so we started reviewing different approaches and came up with good literature and method.

This implementation uses Two Layer networks in which a second layer is the adaptive units which are implemented to the previously mentioned layers. In this, the output is the same as the previous one but in this network, we also have hidden units which do not have any direct effect on the networks environment. The input of the second layer is the output of the first layer, the first layer is the evaluation network and the second layer is the learning layer. As the neural network contains a hidden layer and the output layers, in general, the hidden layers for this structure are the second learning layer and output layers are the original layers. The error propagation to hidden evaluation units and action units depends upon individual networks output. This is a method called error backpropagation which is gradient descent technique for learning nonlinear, differentiable, output function in the hidden unit is very effective. We are working on this approach and reviewing many other papers and resources so that we can solve the above-mentioned issue with the inverted pendulum balancing system

D. Inverted Pendulum using Widrow-Hoff LMS Algorithm

According to [VW88], the force required to stabilize the inverted pendulum system shown in Fig1 at a time instant t is -

$$F = U * \text{sgn}[(k_1 + k_2)x_t + k_2v_{t-1} + (k_3 + k_4)\theta_t - k_4\omega_{t-1}] \quad (7)$$

where : the position of the cart (x), the velocity of the cart (v), the angle of the pendulum (θ), and the angular velocity of the pendulum (ω). The cart and pendulum velocities can be estimated from the instantaneous cart and pendulum velocities which can be obtained from the current states and its history. For the network to balance the pendulum, in some sense it needed to implement Equation 10, estimating the pendulum and cart positions multiplied by the coefficients k_1, k_2, k_3 , and k_4 . Finding a set of weights requires solving a system of M linear equations with N unknowns where M is the number of images and N is the number of weights. Minimum mean-squared-error can be found iteratively by using the Widrow-Hoff LMS algorithm.

Although, the above work used visual images to train the network. It is important to note that training is limited, therefore, machine should be capable enough to perform the task independent of training, responding correctly to almost all situations.

REFERENCES

- [MC68] Donald Michie and R. A. Chambers. "BOXES: An Experiment in Adaptive Control". In: (1968).
- [CL87] K. C. Cheok and N. K. Loh. "A Ball-Balancing Demonstration of Optimal and Disturbance-Accommodating Control". In: *IEEE Control Systems Magazine* 7.1 (1987). DOI: 10.1109/MCS.1987.1105235.
- [VW88] V.V.Tolat and Bernard Widrow. "An adaptive 'broom balancer' with visual inputs". In: *IEEE 1988 International Conference on Neural Networks* (1988). DOI: 10.1109/ICNN.1988.23982.
- [And89] Charles W Anderson. "Learning to control an inverted pendulum using neural networks". In: *IEEE Control Systems Magazine* 9.3 (1989), pp. 31–37. DOI: 10.1109/37.24809.
- [Mic] University of Michigan. *Inverted Pendulum: System Modeling*. URL: <http://ctms.engin.umich.edu/CTMS/index.php?example=InvertedPendulum§ion=SystemModeling>.