

# Machine Learning for Control

by

Rodrigo Queiro (DOW)

Fourth-year undergraduate project in  
Group F, 2010/2011

I hereby declare that, except where specifically indicated, the work submitted herein is my own original work.

Signed: \_\_\_\_\_ Date: \_\_\_\_\_

## Technical Abstract

technical abstract...

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Reinforced Model Learnt Control</b>	<b>4</b>
2.1	Practical Concerns . . . . .	5
2.1.1	GPR Implementation . . . . .	5
2.1.2	Choice of State Vector . . . . .	6
2.1.3	Controller Form . . . . .	7
<b>3</b>	<b>Apparatus and Experimental Results</b>	<b>8</b>
<b>4</b>	<b>Results and Discussion</b>	<b>8</b>
<b>5</b>	<b>Conclusions</b>	<b>8</b>
<b>A</b>	<b>Extra Stuff</b>	<b>10</b>

# 1 Introduction

Balancing a unicycle is a very challenging task for a human rider. Many attempts have been made to achieve this task, using a variety of models for the action of the rider. Some represent the rider as a flywheel or pendulum in the coronal plane, allowing direct compensation of falling to the side [1,2], as in Figure 1(a). Other use a more realistic (and challenging) model of a flywheel in the horizontal plane [3,4], as in Figures 1(b), but none of these have reliably balanced a real unicycle.

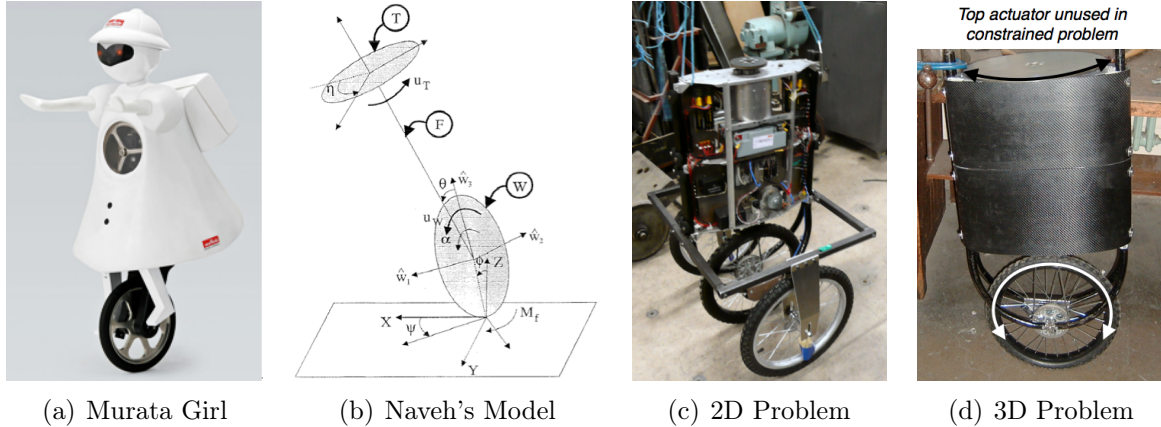


Figure 1: Different balance problems

In 2004/2005 Mellors and Lamb [5,6] built a robotic unicycle, shown in Figure 1(d), intending to design a controller to balance it. However, they were only able to complete the construction of the unicycle. In 2007/2008, D’Souza-Mathew resumed work, replacing a wheel sensor and attempting to design a controller. He simplified the problem by removing the ability to fall to the side, reducing it to 2D dynamic control: the inverted pendulum (from now on referred to as the 2D system). This is shown in Figure 1(c). He was unable to balance the unicycle due to hardware problems.

Next, in 2008/2009 Forster analysed the dynamics of both the 2D problem and the unrestricted 3D unicycle [7]. Again, hardware problems prevented him from balancing the 2D system, and although he designed a controller for the 3D unicycle, it was not even tested in simulation. Given the simplicity of his approach compared to those of Vos and Naveh [3,4], it appears unlikely to work.

One thing all these approaches have in common is that their first step is a series of simplifying assumptions about the dynamic system. They ignore the non-linearities like motor dead-zones and wheel friction that are present in any real-world system, and attempt to design a controller to stabilise the idealised system. In many cases, this approach is very successful. However, D’Souza-Mathew and Forster found that their model was invalid since the unicycle’s motor drive didn’t react faster enough. Vos and Naveh had to use complex, approximate techniques to model the unicycle.

An alternative “intelligent” approach to control involves learning the dynamics of the system directly, instead of relying on assumptions and mechanical analysis. Various methods for this have been used, but many require prohibitively large amounts of data from the system. One method due to Rasmussen and Deisenroth, known here as Reinforced Model Learnt Control (RMLC), achieves unprecedented data efficiency, and has been successfully used to stabilise a computer simulation of a 3D unicycle [8, 9].

In 2009/2010, McHutchon successfully applied RMLC to the 2D system [10]. However, since he had to make significant changes to the unicycle hardware and software to achieve this, he did not have time to attempt to balance the 3D unicycle.

The principle objective of this project is to apply RMLC to the unrestricted 3D unicycle. This includes the solution of problems identified by McHutchon, as well as other problems identified during the project.

## 2 Reinforced Model Learnt Control

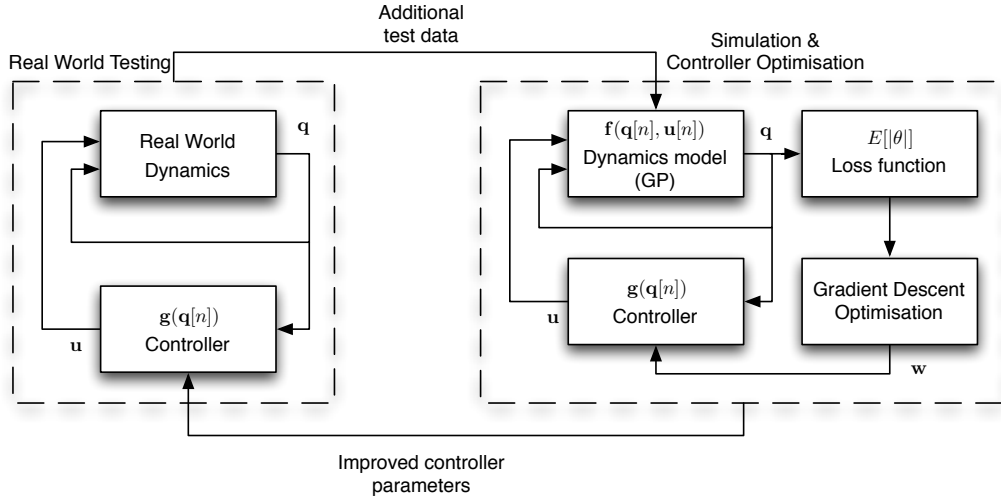


Figure 2: Reinforced Model Learnt Control

The main technique in this project is Reinforced Model Learnt Control, diagrammed in Figure 2. At its core, it assumes that the system (in this case, the unicycle) can be modelled in discrete time as:

$$\mathbf{q}[n + 1] = \mathbf{f}(\mathbf{q}[n], \mathbf{u}[n])$$

In this equation,  $\mathbf{q}[n]$  is the state of the system at time  $n$ , and  $\mathbf{u}[n]$  is control input at time  $n$ . In the case of the unicycle,  $\mathbf{q}$  consists of angles and angular velocities of the components of the unicycle, and the position of the unicycle.  $\mathbf{u}$  consists of the commands sent to the wheel and flywheel motors.

This function,  $\mathbf{f}$ , is modelled as a Gaussian Process (GP). By using Gaussian Process Regression (GPR), we can estimate any continuous function from sampled inputs and outputs. For a description of the mechanics of GPR, refer to [11]. When the unicycle runs, we get a series of states and control inputs that can be converted to samples of  $\mathbf{f}$ , and this allows us to use GPR to estimate  $\mathbf{f}$  at any point. This estimated  $\mathbf{f}$  is referred to as the **dynamics model**.

By successively applying  $\mathbf{f}$  to an initial distribution of possible starting states, we can estimate, with confidence bounds, a distribution of states over some finite horizon. This is referred to as **simulation** of the system. Then, a **loss function** is applied to the state distributions—this might find, for example, the expected distance between the top of the unicycle and the upright position. Summing these losses over the horizon gives a numerical score that rates how well the dynamics model believes a given controller will balance the unicycle. This loss score penalises uncertainty as well as falling.

The gradient of the loss with respect to the controller parameters can be calculated, and this allows standard gradient descent optimisation methods to be used to find a locally optimal controller (for the estimated dynamics model). This process is shown as the right-hand box, “Simulation & Controller Optimisation”, in Figure 2, and is referred to as **training** a controller.

Once a optimal controller has been trained on the simulated system, a **rollout** is performed on the real system (“Real World Testing” in Figure 2). This generates a log of states and control inputs, which can be converted into more samples of  $\mathbf{f}$ , improving the quality of the dynamics model and allowing a better controller to be trained. This process is repeated iteratively until the dynamics model is sufficiently accurate that the trained controllers perform well on the real system.

## 2.1 Practical Concerns

There are many different decisions to make when implementing the RMLC strategy, which are detailed in this section. Fortunately, due to previous work on stabilising both the real 2D system and the simulation of the 3D unicycle, we had a lot of information on which choices can work well in these situations, and which are most important.

### 2.1.1 GPR Implementation

The GPR system has many configurable parameters, but fortunately the form used previously had proven very robust. This project uses a zero-mean GP with a squared exponential covariance function, with automatic relevance detection. This has the form:

$$k(\mathbf{x}, \mathbf{x}') = \alpha^2 \exp \left( - \sum_{d=1}^D \frac{(x_d - x'_d)^2}{2l_d^2} \right)$$

This expression contains hyperparameters for the signal variance  $\alpha^2$  and the length scales  $l_d$ . An additional hyperparameter involved in the regression is the noise variance,  $\sigma_\epsilon^2$ . These parameters are chosen to best fit the data without overfitting with the maximum likelihood (ML) method [11]. The optimal values of these hyperparameters are very useful for interpreting how accurate the dynamics model is: a high SNR  $\frac{\alpha}{\sigma_\epsilon}$  suggests the model can predict very well. Furthermore, automatic relevance detection is provided by the length scales - if a variable is not useful for predicting, the ML length scale will tend to inf.

### 2.1.2 Choice of State Vector

The state vector  $\mathbf{q}[n]$  should be chosen to ensure that the future states are a function only of the current state, and current and future control inputs. In other words, the states should form a Markov Chain:

$$P(\mathbf{q}[n+1]|\mathbf{q}[i] \text{ for } i = 1, \dots, n) = P(\mathbf{q}[n+1]|\mathbf{q}[n])$$

Forster's analysis suggested the following state to be suitable, which was found by Rasmussen to be sufficient to model the simulation of the 3D unicycle.

$$\mathbf{q}[n] = \begin{bmatrix} \dot{\theta} & \text{roll angular velocity} \\ \dot{\phi} & \text{yaw angular velocity} \\ \dot{\psi}_w & \text{wheel angular velocity} \\ \dot{\psi}_f & \text{pitch angular velocity} \\ \dot{\psi}_t & \text{flywheel (turntable) angular velocity} \\ x_c & \text{position of target in unicycle's reference frame} \\ y_c & \\ \theta & \text{roll angle} \\ \psi_f & \text{pitch angle} \end{bmatrix}$$

However, this ignores the presence of unobserved states in the system like delays, backlash in the gears, etc. To help the dynamics model deal with these problems, we tried giving it access to the previous state and control input, in effect modelling it as a 2<sup>nd</sup> order Markov chain:

$$\mathbf{q}_2[n] = \begin{bmatrix} \mathbf{q}[n-1] \\ \mathbf{u}[n-1] \\ \mathbf{q}[n] \end{bmatrix}$$

This significantly improved the accuracy of the dynamics model, which in turn suggests that unobserved states are significant in the behaviour of the real unicycle.

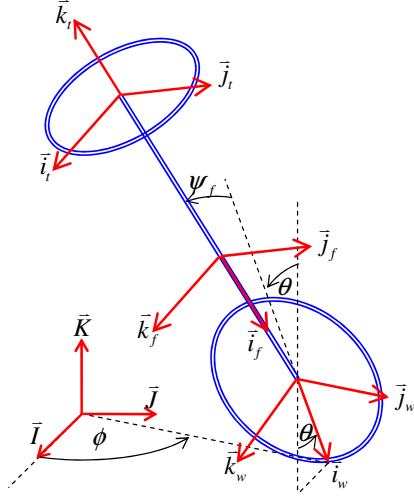


Figure 3: Diagram showing Euler angles for the rotation of the unicycle (from [7])

### 2.1.3 Controller Form

The most basic form of controller is a linear controller,  $\mathbf{g}(\mathbf{q}[n]) = \mathbf{W}\mathbf{q}[n] + \mathbf{p}$ , where  $\mathbf{W}$  is a matrix of weights and  $\mathbf{p}$  is a vector of offsets. This form is capable of stabilising the ideal inverted pendulum, and indeed proved sufficient to stabilise the 2D system.

However, the linear controller cannot generate the correct turning command as shown in Figure 4—this is equivalent to the XOR problem, and can be solved by using a quadratic controller. This takes the form:

$$g_i(\mathbf{q}[n]) = p_i + \sum_{j=1}^D w_{i,j} q_j[n] + \sum_{j=1}^D \sum_{k=j}^D h_{i,j,k} q_j[n] q_k[n]$$

It was found that the controllers performed significantly better when using  $\mathbf{q}_2[n]$  as input, instead of  $\mathbf{q}[n]$ . To understand this, consider the effect with a linear policy: the controller for the augmented state is equivalent to a combination of a two-tap FIR filter and a first-order IIR filter on the original state:

$$\begin{aligned} \mathbf{u}[n] &= \mathbf{g}(\mathbf{q}_2[n]) = \mathbf{W}\mathbf{q}_2[n] \\ &= \mathbf{W}_1\mathbf{q}[n-1] + \mathbf{W}_2\mathbf{u}[n-1] + \mathbf{W}_3\mathbf{q}[n] \end{aligned}$$

This allows the RMLC system to create basic low or high-pass filters in the controller, and this additional freedom improved the subjective quality of the controllers trained.

- GPR implementation (choice of covariance function, sparse approximations, etc.)
- contents of the state vector

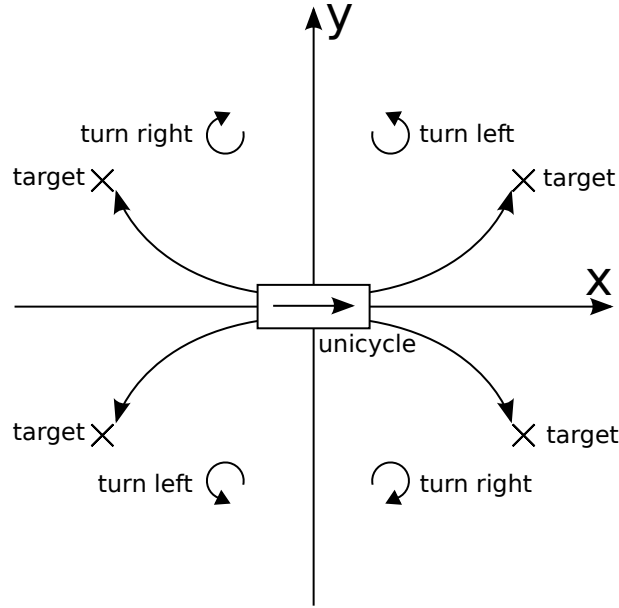


Figure 4: Correct turning command in unicycle-centred coordinates

- form of the controller (e.g. linear:  $\mathbf{g}(\mathbf{q}[n]) = \mathbf{w}^T \mathbf{q}[n]$ )
- loss function
- timestep for conversion to discrete time

### 3 Apparatus and Experimental Results

### 4 Results and Discussion

### 5 Conclusions

## References

- [1] D. Zenkov et al. The Lyapunov-Malkin theorem and stabilization of the unicycle with rider. *Systems and Control Letters*, 46:293–302, 2002.
- [2] Murata Manufacturing Co. Development of the unicycle-riding robot: Murata girl. [http://www.murata.com/new/news\\_release/2008/0923.html](http://www.murata.com/new/news_release/2008/0923.html), 2008.
- [3] D. Vos. Nonlinear control of an autonomous unicycle robot: Practical issues. MIT PhD Thesis, 1992.
- [4] Y. Naveh et al. Nonlinear modelling and control of a unicycle. *Dynamics and Control*, 9:279–296, 1999.



- [5] M. Mellors. Robotic unicycle: Mechanics & control. CUED Master's Project, 2005.
- [6] A. Lamb. Robotic unicycle: Mechanics & control. CUED Master's Project, 2005.
- [7] D. Forster. Robotic unicycle. CUED Master's Project, 2009.
- [8] Carl Edward Rasmussen and Marc Peter Deisenroth. Recent advances in reinforcement learning. chapter Probabilistic Inference for Fast Learning in Control, pages 229–242. Springer-Verlag, Berlin, Heidelberg, 2008.
- [9] Marc P. Deisenroth and Carl E. Rasmussen. PILCO: A Model-Based and Data-Efficient Approach to Policy Search. In L. Getoor and T. Scheffer, editors, *Proceedings of the 28th International Conference on Machine Learning*, Bellevue, WA, USA, June 2011.
- [10] A. McHutchon. Machine learning for control. CUED Master's Project, 2010.
- [11] Carl Edward Rasmussen. Gaussian processes for machine learning. MIT Press, 2006.

## A Extra Stuff