
REINFORCEMENT LEARNING FOR CONTROL AND MULTIPLAYER GAMES

IIB PROJECT INVESTIGATING THE APPLICATION OF GOOGLE DEEPMIND'S
ALPHAZERO ALGORITHM TO CLASSICAL CONTROL PROBLEMS



**UNIVERSITY OF
CAMBRIDGE**
DEPARTMENT OF ENGINEERING

AUTHOR **ALEX DARCH**

SUPERVISOR **DR. GLENN VINNICOMBE**

ASSESSOR **DR. IOANNIS LESTAS**

*St John's College
Cambridge
May 12, 2019*

0 Contents

1	Introduction	3
1.1	Controlling Dynamical Systems	3
1.2	Control Theory	3
1.3	Reinforcement Learning	4
1.4	AlphaZero	5
1.4.1	Self Play and The Neural Network	6
1.4.2	Monte Carlo Tree Search	7
1.4.3	Policy Iteration	8
2	Theory and Methods	9
2.1	The Inverted Pendulum (IP)	9
2.1.1	Dynamics	9
2.1.2	Cost and Value Function	11
2.1.3	State Representations	11
2.1.4	Discrete vs Continuous Time	13
2.2	Self Play and Adversaries	14
2.2.1	Cost and Value Functions of the Adversary	14
2.2.2	Choice of Action	14
2.2.3	Episode Execution	15
2.3	Neural Network	15
2.3.1	Loss Functions and Pareto	15
2.3.2	Architectures	15
2.4	MCTS	15
2.4.1	State and Player Representation	16
2.4.2	Terminal States and Suicide***	16

2.4.3	Modified UCB	16
2.5	Player and Adversary Evaluation	16
2.5.1	Elo Scoring	16
References		17
A Appendices		19
A.1	Inverted Pendulum Dynamics Derivation	19
A.1.1	Swing Up Control	20
A.2	Propagation of Quantisation Error	20

1 Introduction

1.1 Controlling Dynamical Systems

Dynamical systems have long been of great interest to engineers and scientists due to their ability to handily describe real world phenomena. They describe how a system changes through geometrical space with time and, in principle, their trajectories can be predicted solely from their initial state. In recent years this has become relatively easy, even for non-linear systems, by using numerical methods coupled with increases in computing power. However, often a precise solution is of less import than a prediction of the long term qualitative behaviour of the system. Much effort has been made to find methods to describe this long-term behaviour such as those made by Lyapunov in stability theory.

A corollary of the importance of dynamical systems is that influencing their behaviour is particularly useful. Highly non-linear models are difficult to control optimally and robustly, and the few mathematical techniques developed to deal with these can only handle very small subcategories of problems [1, 2]. The study of more general techniques to solve complex control problems has recently come to the forefront of the field with the advent of machine learning techniques such as reinforcement and deep learning. Many of these are not aimed specifically at control problems and are often designed to play games. This project looks at adapting one such algorithm - AlphaZero - from playing board games to solving general control problems such as the control of an aircraft in flight or the control of under-actuated robotics by starting with a simple inverted pendulum.

1.2 Control Theory

The first attempts at controlling dynamical systems came from classical control theory, which consisted of a series of 'cut-and-try' techniques based largely on adjusting the gains of PID controllers and lead/lag compensators until satisfactory closed loop dynamics were achieved [4].

It wasn't until the unification of the calculus of variations, classical control, random process theory and linear/non-linear programming by bellman in the 1950's [4] that truly optimal control was discovered. Optimal control consists of finding a control law for a

specified optimality criterion and can be thought of as a non-linear feedback law, $u(t) = -K(t)x(t)$ based on the initial state, $x(0)$. In discrete time, an optimal control law can be found via dynamic programming (DP). DP is a systematic procedure for transforming an optimisation over a sequence of h inputs into h minimisations over 1 input (but for all states).

$$V(x_k, k) = \min_{u_{k:h-1}} \left(\sum_{i=k}^{h-1} c(x_i, u_i) + J_h(x_h) \right) \quad (1.1)$$

$$= \min_{u_k} \left(c(x_k, u_k) + V(x_{k+1}, k+1) \right) \quad (1.2)$$

$$u_k^* = \underset{u_k}{\operatorname{argmin}} \left(c(x_k, u_k) + V(x_{k+1}, k+1) \right) \quad (1.3)$$

The dynamic programming equations above, where $c(x_i, u_i)$ is the cost as a function of the state and input at time i , $J_h(x_h)$ is the terminal cost, and $V(x_k, k)$ is the value function.

This enables a backwards recursion to find a sequence of value functions. This can be solved for the linear case with quadratic costs analytically (LQR) or, if the system is non-linear, via gradient descent. However, over a finite horizon this is essentially open-loop control. Optimal control of this form has been used to control complex dynamical systems such as spaceflight and aileron folding on aircraft [5, 6]. A closed loop extension to this is Model Predictive Control (MPC), which employs a receding horizon rather than a finite or infinite horizon. MPC can therefore easily deal with plant disturbances and uncertainties, constraints, indefinite horizons and can also be extended to get a control law for non-linear systems. MPC has recently been shown to work in trajectory control for interplanetary rovers [7]. Optimal control is limited in requiring sophisticated models of the environment/plant and generally struggle with highly non-linear models - state of the art is currently linearisation about the predicted trajectory. Furthermore, it is only feasible to 'grid' up to 5/6 dimensions in discrete cases [2].

1.3 Reinforcement Learning

Two further generalisations to dynamical systems and optimal control as defined in eq. (1.1) are stochastic dynamics and indefinite horizons (i.e. episodic tasks). This discrete time stochastic control process is known as a Markov Decision Process (MDP). In MDPs the cost function is often written as a reward function and, due to the indefinite nature of the process, the value function for the next step is discounted (where $\lambda \approx 0.9$ typically):

$$V(x_k) = \max_{u_k} \left(\sum_{i=k}^{\infty} \lambda^{i-k} r(x_i, u_i) \right) \quad (1.4)$$

$$= \max_{u_k} \mathbb{E} \left[r(x_k, u_k) + \lambda V(x_{k+1}) \right] \quad (1.5)$$

$$u_k^* = \operatorname{argmax}_{u_k} \mathbb{E} \left[r(x_k, u_k) + \lambda V(x_{k+1}) \right] \quad (1.6)$$

Reinforcement learning (RL) aims to learn the optimal policy, $\pi^*(x_k)$ ($= u^*(x_k)$ in control) of an MDP. This differs from optimal control in its inherent stochastic nature and therefore can lead to intractable search spaces. A solution to this is to learn from sample trajectories. Algorithms such as Q-Learning, SARSA and DYNA have recently had great success in control applications, for example, their use in controlling mobile robots [9, 10]. Furthermore, the advent of neural networks has led to the extension of these to functional approximations from tabula-rasa methods, making the control highly non-linear dynamical systems possible. Notably, Deepmind’s recent success with training a robot to gently manipulate objects [11], would not be possible to reproduce using classical or modern control techniques due to dimensional problems.

1.4 AlphaZero

AlphaGo Zero is a revolutionary Reinforcement Learning algorithm that achieved super-human performance in the game of go, winning 100–0 against the previously published, champion-defeating AlphaGo. Its successor, AlphaZero, is a generalised version that can achieve superhuman performance in many games. There are two key sub-algorithms that form the basis of their success: Monte-Carlo tree search (MCTS) for policy improvement, and a deep CNN for the neural policy and value network. Policy iteration is then implemented through self-play. AlphaGo Zero and AlphaZero differ only in the latter’s use of a single neural network that is updated iteratively, rather than evaluated against the previous one, and by not taking advantage of the symmetries of the games.

A key feature of AlphaZero is that it only requires the ability to simulate the environment. It does not need to be told how to win, nor does it need an exact model of the system dynamics, $p(s_{t+1}|s_t, a_t)$, as this can be learnt through self-play. Furthermore, the algorithm often ‘discovers’ novel solutions to problems, as shown by ‘move 37’ in a game of Go against the reigning world champion, Lee Sedol. This makes it particularly suitable for learning to control complex dynamical systems where approximate simulations can be made.

1.4.1 Self Play and The Neural Network

Figure 1.1a shows how self-play is performed in AlphaZero. An episode/game, $\{s_1, \dots, s_T\}$ is played and for each state, s_t , a Monte-Carlo Tree Search is performed, guided by the current policy f_θ . The MCTS outputs an improved action-*p.m.f*, $\pi_t = [p(a_1|s_t), p(a_2|s_t), \dots, p(a_N|s_t)]$. The next move is then selected by sampling from π_t . The final player then gets scored (e.g in chess $z \in \{-1, 0, 1\}$ for a loss, draw or win respectively). The game score, z , is then appended to each state-action pair depending on who the current player was on that move to give training examples of the form $(s_t, \pi_t, (-1)^{\mathbb{I}(\text{winner})}z)$.

The neural network, $f_\theta(s)$ takes a board state, s , and outputs a *p.m.f* over all actions and the expected outcome, (\mathbf{p}_θ, v) (fig. 1.1). The networks are initialised with $\theta \sim \mathcal{N}(0, \epsilon)$, ϵ small. The neural network is trained to more closely match the MCTS-informed action-*p.m.f*s, π_t , and the expected outcome (state-value), v_t to z .

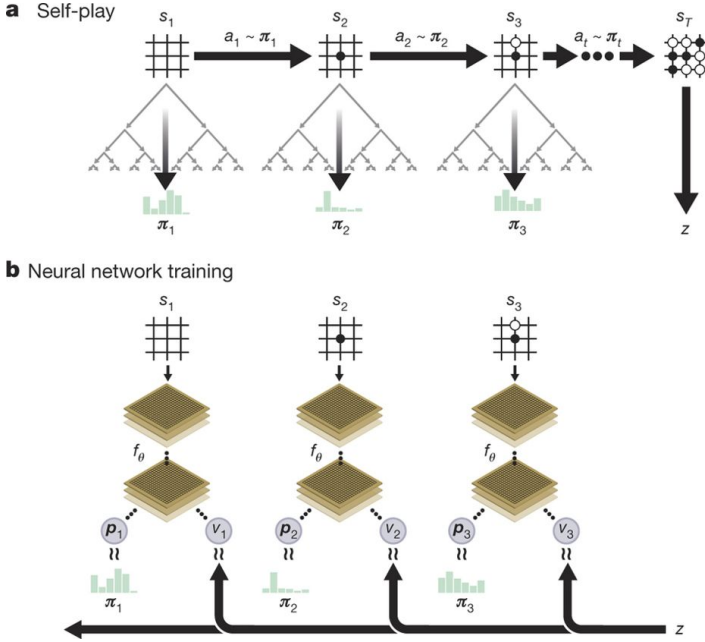


Figure 1.1: A schematic showing how self-play and policy training are performed. Taken from [12].

block and two separate 'heads'. The policy head has softmax activation function, preceded by series of ReLU linear layers and batch normalisation layers. The value head also has this but with a tanh output activation. The input convolutional block consists a single convolutional layer followed by 19-39 residual blocks.

The loss function for the neural network is given by:

$$\mathcal{L} = (z - v)^2 - \pi \cdot \log(\mathbf{p}) + c \|\theta\|^2 \quad (1.7)$$

For chess, the input state consists of an image stack of 119 8x8 planes representing the board state at times $\{t, t-1, \dots, t-8\}$, and planes representing repetitions, colours and castling etc. The output action *p.m.f* is a $8 \times 8 \times 73 = 4672$ vector representing every possible move from each piece, illegal moves are then masked. The output value, $v \in (-1, 1)$. The network itself consists of an input convolutional

1.4.2 Monte Carlo Tree Search

Figure 1.2a: A MCTS is performed at each step of an episode. The state at which the tree search starts then becomes the root state, s_{root} . From the root state, the tree search can move to an edge (s, a) by selecting an action, a . Before selection a prior $P(s, a) = (1 - \epsilon)\mathbf{p}_{\theta, root} + \epsilon\boldsymbol{\eta}$ where $\boldsymbol{\eta} \sim Dir(0.3)$ and $\epsilon = 0.25$ is added to ensure exploration. Each edge stores a prior action-*p.m.f.* output by the policy, $P(s, a) = \mathbf{p}_{\theta}$; a visit count, $N(s, a)$; and an action-value, $Q(s, a)$. Actions are selected by maximising an the action-value plus an upper confidence bound, which encourages exploration. The constant c (~ 1) can be decreased to encourage exploitation.

$$a_{selected} = \underset{\forall a}{argmax} \left\{ Q(s, a) + c \cdot P(s, a) \frac{\sqrt{\sum_b N(s, b)}}{1 + N(s, a)} \right\} \quad (1.8)$$

Figure 1.2b: Once a leaf node ($N(s, a) = 0$) is reached, the neural network is evaluated at that state: $f_{\theta}(s) = (p(s, \cdot), v(s))$. The action *p.m.f.* and state-value are stored for the leaf state.

Figure 1.2c: The action-values, Q , are calculated as the mean of state-values in the subtree below that action. The state-value are then calculated as the mean of all the action-values branching from that state, and so on.

$$Q(s, a) = \frac{1}{N(s, a)} \sum_{s_{t+1}|s_t, a_t} v(s_{t+1}) \quad (1.9)$$

Figure 1.2d: After a set number of simulations (1600), the MCTS-improved action *p.m.f.s*, $\boldsymbol{\pi} = p(\mathbf{a}|s_{root}) \propto N^{1/\tau}$ are returned, where N is the visit count of each move from the root state and τ controls the sparseness of the probability mass function ($\{\tau = 0\} \rightarrow \underset{a}{argmax} \{N(s, a)\}$). For self-play, $\tau = 1$ for the first 30 moves and $\tau \rightarrow 0$ thereafter. For evaluation $\tau \rightarrow 0 \forall t$. If $\tau \rightarrow 0$ then $\boldsymbol{\pi}$ becomes one-hot and the loss function of the neural network makes sense as a prediction of a categorical distribution.

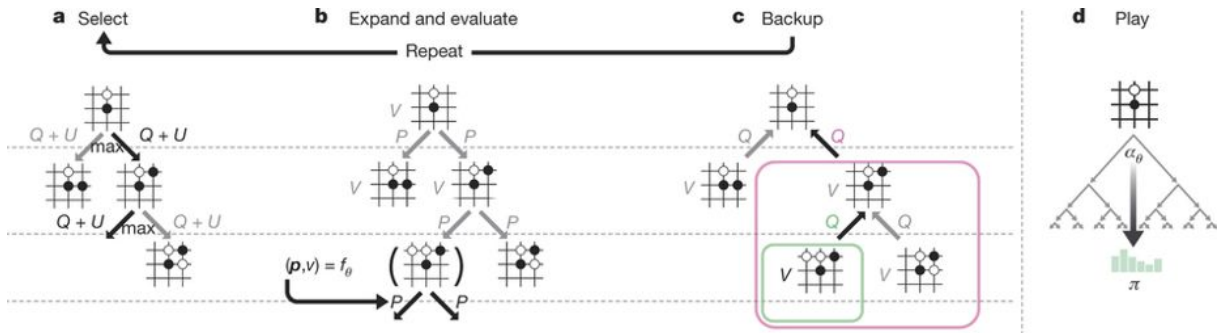


Figure 1.2: A schematic outlining the steps involved in a monte-carlo tree search. Taken from [12].

1.4.3 Policy Iteration

AlphaZero uses a single neural network that continually updates, whether it is better or worse than the previous network. Whereas AlphaGo Zero games are generated using the best player from all previous iterations and then only replaced if the new player wins $> 50\%$ of evaluation games played.

Evaluation is done by playing 400+ greedy ($\tau \rightarrow 0$) games of the current best neural network against the challenging neural network. The networks are then ranked based on an elo scoring system.

The application of AlphaZero to dynamical systems a number of potential benefits above traditional optimal control or reinforcement learning:

Robust Disturbance Handling - The Adversary in AlphaZero is effectively worst case scenario disturbance modelling. By incorporating this into the training process, the controller should be able to cope with situations well outside normal operating conditions.

Non-Linear Systems - Neural networks are universal function approximators and, hence with the correct architecture, therefore should be able to model any system.

General Framework - AlphaZero is a general algorithm that should be able to be applied to many different systems with minimal changes and still model the control well.

This project investigates these.

2 Theory and Methods

Explain the assumptions behind the theoretical development you are using and the application of the theory to your particular problem. Any heavy algebra or details of computing work should go into an appendix. This section should describe the running of the experiment or experiments and what equipment was used, but should not be a blow by blow account of your work. Experimental accuracy could be discussed here.

2.1 The Inverted Pendulum (IP)

The Inverted Pendulum is an inherently unstable system with highly nonlinear dynamics and is under-actuated.

2.1.1 Dynamics

The full state space equations for the inverted pendulum as defined in fig. 2.1 are given by:

$$\begin{bmatrix} \dot{x} \\ \ddot{x} \\ \dot{\theta} \\ \ddot{\theta} \end{bmatrix} = \begin{bmatrix} \dot{x} \\ \frac{\left(\frac{2M-m}{m}F_2 - F_1\right)\cos\theta + g(M+m)\sin\theta - mL\dot{\theta}^2\sin\theta\cos\theta}{(M+m\sin^2\theta)} \\ \dot{\theta} \\ \frac{F_1 + F_2\cos(2\theta) + m\sin\theta(L\dot{\theta}^2 - g\cos\theta)}{L(M+m\sin^2\theta)} \end{bmatrix} \quad (2.1)$$

Ignoring second order terms and linearising about $\mathbf{x}_e = [x_e, \dot{x}_e, \theta_e, \dot{\theta}_e]^T = [0, 0, 0, 0]^T$:

$$\begin{bmatrix} \dot{x} \\ \ddot{x} \\ \dot{\theta} \\ \ddot{\theta} \end{bmatrix} = \begin{bmatrix} \dot{x} \\ \frac{\frac{2M-m}{m}F_1 - F_2 + g(M+m)\theta}{M} \\ \dot{\theta} \\ \frac{F_1 + F_2 - gm\theta}{lM} \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & g\frac{M+m}{M} & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & -\frac{mg}{lM} & 0 \end{bmatrix} \begin{bmatrix} x \\ \dot{x} \\ \theta \\ \dot{\theta} \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ -\frac{1}{M} & \frac{2M-m}{Mm} \\ 0 & 0 \\ \frac{1}{lM} & \frac{1}{lM} \end{bmatrix} \begin{bmatrix} F_1 \\ F_2 \end{bmatrix} \quad (2.2)$$

Which, as expected, is unstable since $\det(\lambda I - A) = 0 \implies \lambda^2(\lambda^2 + \frac{mg}{lM}) = 0$. Note, for small angles the natural frequency of a non-inverted pendulum is $\omega_n = \sqrt{\frac{mg}{lM}} = \sqrt{\frac{0.1 \times 9.81}{0.5 \times 1}} \approx 1.40 \text{ rad/s}$. Therefore, the time constant for the system is $\tau \approx 0.70 \text{ s}$.

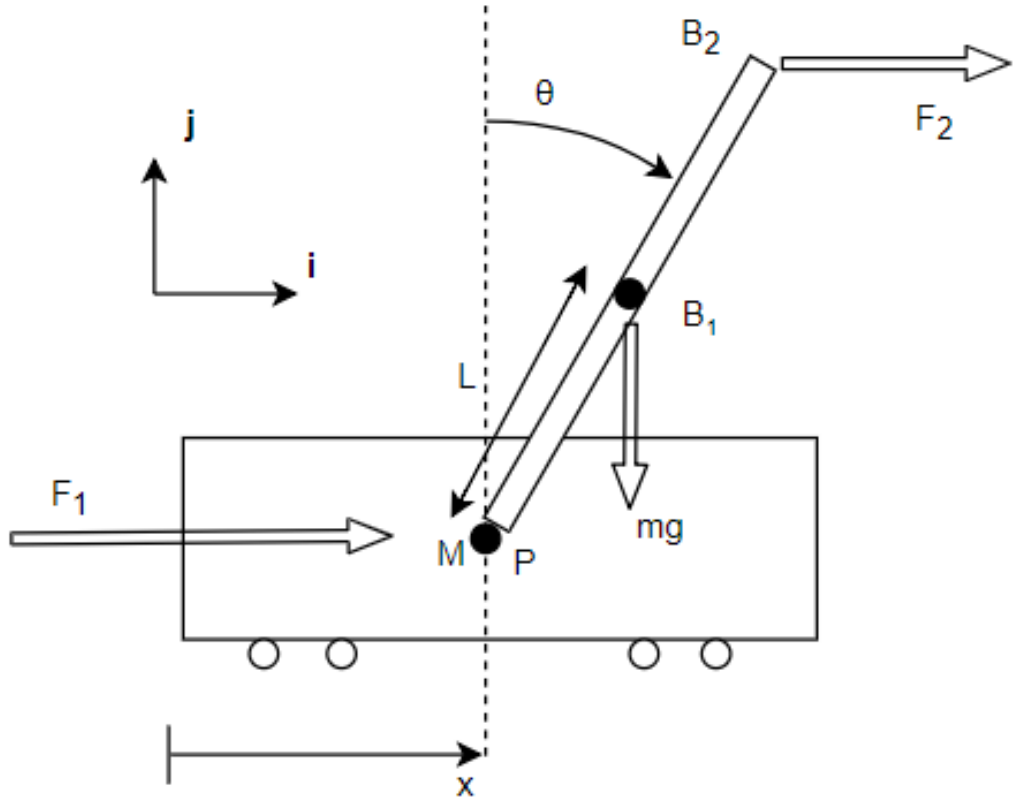


Figure 2.1: A free-body diagram of the inverted pendulum system. For the OpenAI IP the system is in discrete time with a time-step of $\tau = 0.02s$. The other constants are $l = 0.5m$, $m = 0.1kg$, $M = 1kg$, $F = \pm 10N$, $x_{max} = \pm 2.4m$, $\theta_{max} = \pm 12^\circ$.

OpenAI's gym is a python package that supplies an inverted pendulum environment built-in. This environment was wrapped to use the dynamics above and other extra functionality, whilst providing a rendering function shown in fig. 2.2.

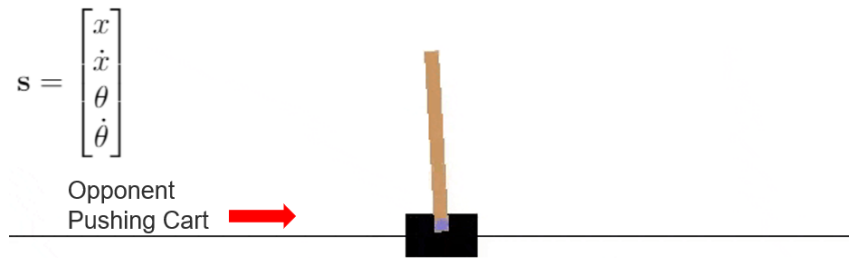


Figure 2.2: The OpenAI gym CartPole environment. The classical state representation is shown in the top left. Actions by the player and the adversary are taken as an impulse to the left or right as defined in fig. 2.1.

2.1.2 Cost and Value Function

For each step/impulse, the 2D state is calculated and a cost, is calculated as:

$$c(x_t, u_t) = -\frac{1}{\sum_i w_i} \mathbf{w} \cdot \left[\left(\frac{x_t}{x_{max}} \right)^2, \left(\frac{\dot{x}_t}{\dot{x}_{max}} \right)^2, \left(\frac{\theta_t}{\theta_{max}} \right)^2, \left(\frac{\dot{\theta}_t}{\dot{\theta}_{max}} \right)^2 \right]^T \quad (2.3)$$

Where $\mathbf{w}^T = [w_1, w_2, w_3, w_4] = [0.25, 0.1, 0.7, 1]$ and $0 \geq c(x_t, u_t) \geq -1$. The weights, \mathbf{w} , were chosen through empirical measurement of the the importance of each state ***. Weighting on the inputs was set to zero, as there are only two inputs for this problem, thus the cost can be written as $c(x_t)$. The max values can be approximated experimentally (note, $x_{max} = 2.4$ and $\theta_{max} = 12^\circ$ are given constraints):

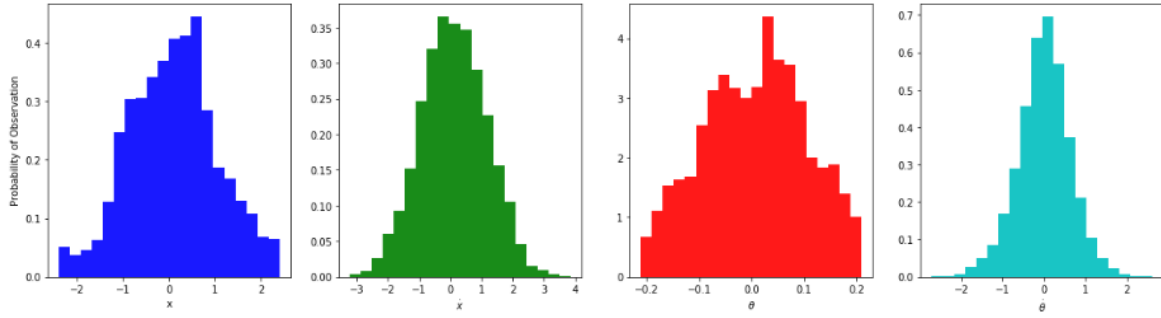


Figure 2.3: Histograms of typical state values. The frequencies greatly depend on the quality of the controller, with better controllers giving much narrower distributions. However, these are typical for a controller of medium efficacy over many episodes where the starting state is randomised (possibly to an uncontrollable state).

Suitable estimates for the the values of \dot{x}_{max} and $\dot{\theta}_{max}$ are thus ≈ 3 and 2 respectively.

The value function is computed after an episode has completed as the discounted future losses at each state with the constraint that $\gamma^k < \frac{1}{20}$, where $\frac{1}{20}$ was chosen as it is a standard factor for insignificance. Since steps_beyonds_done (= k) must be defined in the CartPoleWrapper class, this is a constant, and therefore γ is calculated as $\gamma < \frac{1}{20^{\frac{1}{k}}}$. The discounted values are calculated using a geometric series:

$$v_0 = \frac{\sum_{\tau=0}^k \gamma^\tau c(x_\tau)}{\sum_{\tau=0}^k \gamma^\tau}, \quad \text{where } \gamma^k < \frac{1}{20} \quad (2.4)$$

Where for simplicity of notation, $v_0 = v(t)$, the state value at step t.

2.1.3 State Representations

The state can be represented in a number of ways, most simply this would just be feeding $\mathbf{x} = [x, \dot{x}, \theta, \dot{\theta}]$ into the neural network. This has a number of advantages such as lower

computational cost, greater numerical accuracy (if the process is fully observable) and simpler implementation. Conversely, following Silver et. al, a 2D representation can be used. There are several possibilities for this, all of which first require binning \mathbf{x} :

(1) A matrix stack of \mathbf{x} vs $\dot{\mathbf{x}}$ and θ vs $\dot{\theta}$, both of which would only have one non-zero entry. This scales as b^n where b = number of bins and n = number of states.

(2) A matrix stack of x_t vs x_{t-1} for all states. Similarly this scales as b^n , however the derivative states do not need to be plotted as these can be inferred. This has the advantage that, if the derivatives are not observable, we can build them into the 2D representation, however, if they are observable then this is less accurate than (1).

(3) A matrix of discounted previous states forming a motion history image. This is the formulation used, and shown in fig. 2.4. algorithm 1 shows the implementation details.

A 2D representation like this allows us to use a convolutional neural network, which has the benefit of various transformation invariances - these are particularly useful for the inverted pendulum since it is symmetric.

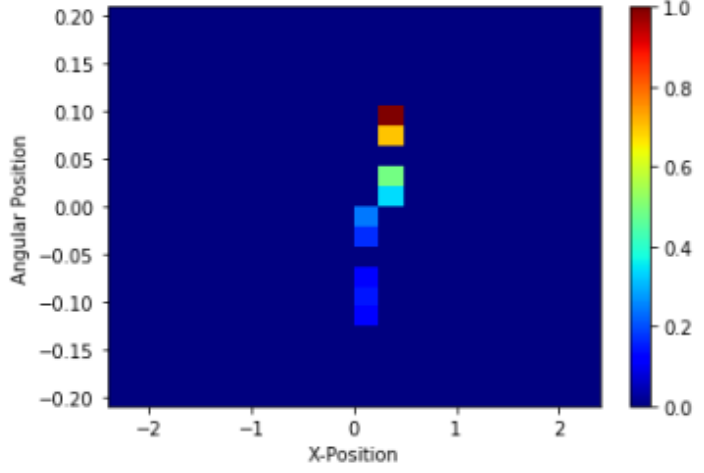


Figure 2.4: An example of a 2D state representation where there are 20 bins and 17 random actions have been taken.

Algorithm 1 Create 2D State

```

1: function GETSTATE2D( $\hat{\mathbf{x}}_{t-1}^{(2D)}$ , binEdges, nBins)
2:    $\hat{\mathbf{x}} \leftarrow getNormedState()$ 
3:   for all  $x_i \in \hat{\mathbf{x}}$  do
4:      $x_i \leftarrow argmin|binEdges - x_i|$  ▷ get the index of the nearest bin edge.
5:   end for
6:    $HistEdges \leftarrow linspace(-0.5, nBins - 0.5, nBins + 1)$  ▷ centre by shifting -0.5
7:    $\hat{\mathbf{x}}_t^{(2D)} \leftarrow histogram2d(x, \theta, bins = (HistEdges, HistEdges))$  ▷ Inbuilt function
8:    $\hat{\mathbf{x}}_{t-1}^{(2D)} [\hat{\mathbf{x}}_{t-1}^{(2D)} < \lambda^{-(T+1)}] \leftarrow 0$  ▷ Only keep  $\hat{\mathbf{x}}_{t-1}^{(2D)}$  from  $t < T$ 
9:   return  $\hat{\mathbf{x}}_t^{(2D)} + \lambda \hat{\mathbf{x}}_{t-1}^{(2D)}$ 
10: end function

```

The state space model for the quantisation of the linearised inverted pendulum (valid for small time steps) can be modelled as:

$$\mathbf{x}_t^{(2D)} = C\mathbf{x}_t + \mathbf{V}_t \quad \mathbf{V}_t \sim WN(0, \sigma_v^2 I) \quad (2.5)$$

$$\mathbf{x}_t = A\mathbf{x}_{t-1} + B\mathbf{u}_t \quad (2.6)$$

The initial mean squared error, and the propagation of the error (when estimated optimally) are given by eqs. (2.7) and (2.8) (derivation details can be found in appendix A.2).

$$\sigma_0^2 = \sigma_v^2 = \frac{\delta x^2}{12} \quad (2.7)$$

$$\sigma_{n+1}^2 = A\sigma_n^2 \left(1 - \frac{A\sigma_n^2}{A\sigma_n^2 + \sigma_v^2} \right) \quad (2.8)$$

If the spectral density of eq. (2.8) is less than one then the error will decay to zero, and the neural network (if acting as an optimal filter) should be able to recover \mathbf{x} without loss. The number of steps needed for this recovery decreases with smaller bin sizes since $\lim_{\delta x \rightarrow 0} \sigma_v^2 = 0$, and σ_v^2 decays with δx^2 .

With limited memory and a non-zero bin size the overall error can be reduced by binning more densely in regions in which we expect to spend more time. Figure 2.3 shows that the state visit frequencies roughly correspond to normal distribution, therefore by transforming the bins with an inverse gaussian c.f.d. a flattened distribution can be obtained with a greater density of bins in the centre (fig. 2.5). This has the additional benefit of allowing finer control where it is needed. For example, if the pole is far from the equilibrium the optimal action more easily determined, and subject to less change with small state variations.

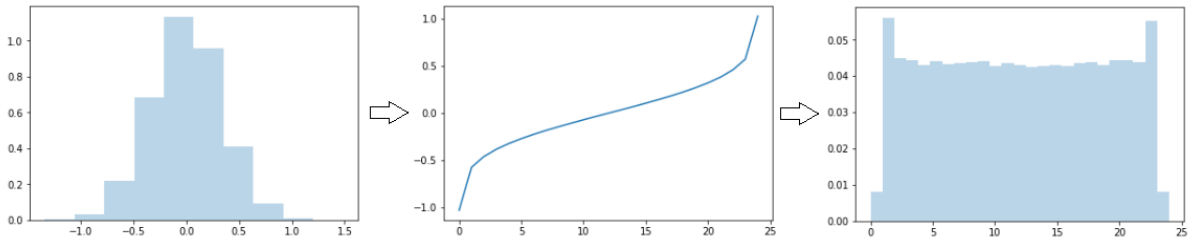


Figure 2.5: Binning of a gaussian distribution with bin edges scaled with an inverse gaussian c.f.d. For this example there are 25 bins.

2.1.4 Discrete vs Continuous Time

For a small enough discrete time step, the simulation will approximate the continuous process very well. Additionally, we can achieve pseudo-continuous actions within the

Work
out
how
small
vs
how
many
steps
(pos-
sibly
with a
second
order
ap-
prox?)

constraint of $u \in \{-1, 1\}$ via pulse-width modulation. A time step of 0.02s is 35x smaller than the time constant of the linearised system ($\tau \approx 0.70s$). The positional change per time step is therefore $\sim \frac{1}{35}\% = 3\%$. Alternatively the average maximum velocities are $3m/s$ and $2rad/s$, thus the maximum positional change per time step is $\dot{x}_{i,max} \times \frac{\tau}{x_{i,max}} = 2.5\%$ and 30% respectively. In practice a 30% change in θ would only occur around θ_{max} and would be uncontrollable, therefore we can take the change to be far less than this. Therefore, we can model the quick succession of actions as a continuous one:

$$F(t) = g(t) * \int_0^\infty u(t)dt \approx g(t) \sum_{n=0}^\infty I\tau\delta(t - n\tau) \quad (2.9)$$

2.2 Self Play and Adversaries

2.2.1 Cost and Value Functions of the Adversary

In AlphaZero, the adversary is working to minimise the value function and the player is working to maximise it. For board games where players are on equal footing a value function representing the expected outcome, $-1 < v < 1$ can be used. This has the advantage of symmetry - when the board is viewed from the other player's perspective the value can be multiplied by -1, and we get the expected outcome for that player.

The adversary for the inverted pendulum has the additional advantage of gravity, making the play un-symmetric. Given that both the state-value and cost are $-1 < v, c < 0$, multiplying by -1 would mean the adversary is maximising a value function $0 < v < 1$. State-values outside these ranges have no meaning. If a single neural network is used, values close to equilibrium may be predicted into the wrong range. Consequently, both the adversary and the player must predict true state-values. This also has the advantage of maintaining an intuitive meaning of the state-value.

2.2.2 Choice of Action

The inverted pendulum system is symmetric in both x and θ allowing for the possibility of taking advantage of this, however as shown with AlphaZero [13], this provides minimal benefit and also hinders generalisation. The adversarial point of action was chosen to be at the top of the pole, acting horizontally (fig. 2.1). Thus ensuring that two distinct policies must be learnt, rather than one being just inverse probabilities of the other. For simplicity $u_{adv} \in \{-1, 1\}$ was chosen.

There are two good ways of representing when it is the adversaries turn for the neural network:

can we? Also proof of uncontrollability?

work out expression for force? Convolution? Just avg with decay term? Ask vinnicombe

Multiply the board representation by -1 such that opponent pieces are negative.

This can only take two players and; A network that outputs both state-values and action $p.m.f.$'s should predict the same state values for both player and adversary, but predict vastly different actions. Therefore, the values will be decoupled from the actions, which was one of the major benefits of using a single neural network. However, a single network is simpler, and this more closely follows AlphaZero's methodology.

Record the player number with each example and use player-labeled examples to train different neural networks. Using a neural network for each player causes half of the examples to be lost as only the relevant player's examples are used to train each network. However, this doesn't suffer from the problems above and can cope with agents with a different number of possible actions more easily.

*** Currently method 2 is being used, but subject to change ***

Note, in the case of the inverted pendulum, the optimal action is the inverse of the worst action. However, this is not a general result, for example, in a system with non-linear and asymmetric dynamics it is possible to have the target perpendicular to the limit of stability, thus for the adversary it is better to push the system into instability rather than away from equilibrium.

2.2.3 Episode Execution

Pseudocode for episode execution following the sections above is shown in algorithm 2.

2.3 Neural Network

2.3.1 Loss Functions and Pareto

do we use mse or not? how do we balance action and value loss if not symmetric? Alpha zero sets them equal, can we use the same principals here?

2.3.2 Architectures

Player vs Adversary Architectures? Combined? GPU, computing power and complexity

2.4 MCTS

outline + pseudocode

Algorithm 2 Execute Episode

```
1: function EXECUTEEPISODE
2:    $example \leftarrow []$ 
3:    $\mathbf{x}, \mathbf{x}^{(2D)}, c \leftarrow resetEpisode()$   $\triangleright$  Set initial  $\mathbf{x}$  randomly and initialise the cost
4:    $player \leftarrow 0$ 
5:   repeat
6:      $\boldsymbol{\pi} \leftarrow getMCTSActionProb(\mathbf{x}, \mathbf{x}^{(2D)}, player)$   $\triangleright$  Perform MCTS Simulations
7:      $example.append((\mathbf{x}^{(2D)}, \boldsymbol{\pi}, c, player))$ 
8:      $u \sim \boldsymbol{\pi}$   $\triangleright$  Sample action
9:      $\mathbf{x}, c \leftarrow step(\mathbf{x}, u)$   $\triangleright$  Take next true episode step
10:     $\mathbf{x}^{(2D)} \leftarrow getState2D(\mathbf{x}, \mathbf{x}^{(2D)})$ 
11:     $player \leftarrow nextPlayer(player)$ 
12:  until  $episodeEnded(\mathbf{x})$ 
13:   $example \leftarrow convertCostsToValues(example)$ 
14:  return  $example$ 
15: end function
```

2.4.1 State and Player Representation

2.4.2 Terminal States and Suicide***

2.4.3 Modified UCB

2.5 Player and Adversary Evaluation

2.5.1 Elo Scoring

A References

- [1] M.C. Smith, I Lestas, *4F2: Robust and Non-Linear Control* Cambridge University Engineering Department, 2019
- [2] G. Vinnicombe, K. Glover, F. Forni, *4F3: Optimal and Predictive Control* Cambridge University Engineering Department, 2019
- [3] S. Singh, *4F7: Statistical Signal Analysis* Cambridge University Engineering Department, 2019
- [4] Arthur E. Bryson Jr, *Optimal Control - 1950 to 1985*. IEEE Control Systems, 0272-1708/95 pg.26-33, 1996.
- [5] I. Michael Ross, Ronald J. Proulx, and Mark Karpenko, *Unscented Optimal Control for Space Flight*. ISSFD S12-5, 2014.
- [6] Zheng Jie Wang, Shijun Guo, Wei Li, *Modeling, Simulation and Optimal Control for an Aircraft of Aileron-less Folding Wing* WSEAS TRANSACTIONS on SYSTEMS and CONTROL, ISSN: 1991-8763, 10:3, 2008
- [7] Giovanni Binet, Rainer Krenn and Alberto Bemporad, *Model Predictive Control Applications for Planetary Rovers*. imtlucca, 2012.
- [8] Richard S. Sutton and Andrew G. Barto, *Reinforcement Learning: An Introduction (2nd Edition)*. The MIT Press, Cambridge, Massachusetts, London, England
- [9] I. Carlucho, M. De Paula, S. Villar, G. Acosta . *Incremental Q-learning strategy for adaptive PID control of mobile robots*. Expert Systems with Applications. 80. 10.1016, 2017
- [10] Yuxi Li, *Deep Reinforcement Learning: An Overview*. CoRR, abs/1810.06339, 2018.

- [11] Sandy H. Huang, Martina Zambelli, Jackie Kay, Murilo F. Martins, Yuval Tassa, Patrick M. Pilarski, Raia Hadsell, *Learning Gentle Object Manipulation with Curiosity-Driven Deep Reinforcement Learning*. arXiv 2019.
- [12] David Silver, Julian Schrittwieser, Karen Simonyan et al, *Mastering the game of Go without human knowledge*. Nature, vol. 550, pg.354–359, 2017.
- [13] David Silver, Thomas Hubert, Julian Schrittwieser et al, *A general reinforcement learning algorithm that masters chess, shogi and Go through self-pla*. Science 362:6419, pg.1140-1144, 2018.

A Appendices

A.1 Inverted Pendulum Dynamics Derivation

We can find the state space equations for the Inverted Pendulum using d'Alembert forces. Firstly we define the distance and velocity vectors to the important points:

$$\begin{aligned}\mathbf{r}_P &= x\mathbf{i} \\ \mathbf{r}_{B_1/P} &= L\sin\theta\mathbf{i} + L\cos\theta\mathbf{j} \\ \mathbf{r}_{B_1} &= (x + L\cos\theta)\mathbf{i} + L\sin\theta\mathbf{j} \\ \dot{\mathbf{r}}_{B_1} &= (\dot{x} + L\dot{\theta}\cos\theta)\mathbf{i} - L\dot{\theta}\sin\theta\mathbf{j}\end{aligned}$$

Linear Momentum, $\boldsymbol{\rho} = \sum_i m_i \dot{\mathbf{r}}_{i/o} = m\dot{\mathbf{r}}_{B_1} + M\dot{\mathbf{r}}_P$:

$$\boldsymbol{\rho} = \begin{bmatrix} (M + m)\dot{x} + mL\dot{\theta}\cos\theta \\ -mL\dot{\theta}\sin\theta \\ 0 \end{bmatrix}$$

Moment of momentum about P, $\mathbf{h}_P = \mathbf{r}_{B_1/P} \times m\dot{\mathbf{r}}_{B_1}$:

$$\begin{aligned}\mathbf{h}_P &= -mL(L\dot{\theta} + \dot{x}\cos\theta)\mathbf{k} \\ \therefore \dot{\mathbf{h}}_P &= -mL(L\ddot{\theta} + \ddot{x}\cos\theta - \dot{x}\dot{\theta}\sin\theta)\mathbf{k}\end{aligned}$$

We can balance moments using $\dot{\mathbf{h}}_P + \dot{\mathbf{r}}_P \times \boldsymbol{\rho} = \mathbf{Q}_e$ and $\mathbf{Q}_e = \mathbf{r}_{B_1/P} \times -mg\mathbf{j} + \mathbf{r}_{B_2/P} \times F_2\mathbf{i}$:

$$\dot{\mathbf{h}}_P + \dot{\mathbf{r}}_P \times \boldsymbol{\rho} = \begin{bmatrix} 0 \\ 0 \\ -mL(\ddot{x}\cos\theta + L\ddot{\theta}) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ -L(mg\sin\theta + 2F_2\cos\theta) \end{bmatrix} = \mathbf{Q}_e$$

And also balance linear momentum using $\mathbf{F}_e = \dot{\boldsymbol{\rho}}$:

$$\dot{\boldsymbol{\rho}} = \begin{bmatrix} (m + M)\ddot{x} + mL(\ddot{\theta}\cos\theta - \dot{\theta}^2\sin\theta) \\ -mL(\ddot{\theta}\sin\theta + \dot{\theta}^2\cos\theta) \\ 0 \end{bmatrix} = \begin{bmatrix} F_1 + F_2 \\ R - mg \\ 0 \end{bmatrix} = \mathbf{F}_e$$

Finally we can write the system dynamics in terms of $\ddot{\theta}$ and \ddot{x} :

$$\ddot{\theta}(M + m\sin^2\theta)L = \left(\frac{2M-m}{m}F_2 - F_1\right)\cos\theta + g(M+m)\sin\theta - mL\dot{\theta}^2\sin\theta\cos\theta \quad (\text{A.1})$$

$$\ddot{x}(M + m\sin^2\theta) = F_1 + F_2\cos(2\theta) + m\sin\theta(L\dot{\theta}^2 - g\cos\theta) \quad (\text{A.2})$$

Simplifying this for our problem by substituting in constants, we can write the full state space equation:

$$\begin{bmatrix} \dot{x} \\ \ddot{x} \\ \dot{\theta} \\ \ddot{\theta} \end{bmatrix} = \begin{bmatrix} \dot{x} \\ \frac{\left(\frac{2M-m}{m}F_2 - F_1\right)\cos\theta + g(M+m)\sin\theta - mL\dot{\theta}^2\sin\theta\cos\theta}{(M+m\sin^2\theta)} \\ \dot{\theta} \\ \frac{F_1 + F_2\cos(2\theta) + m\sin\theta(L\dot{\theta}^2 - g\cos\theta)}{L(M+m\sin^2\theta)} \end{bmatrix} \quad (\text{A.3})$$

It can be proved that the inverted pendulum system is controllable by showing:

$$\text{rank}[\mathbf{B} \ \mathbf{AB} \ \mathbf{A}^2\mathbf{B} \ \mathbf{A}^3\mathbf{B}] = 4 \quad (\text{A.4})$$

Therefore for any initial condition we can reach \mathbf{x}_e in finite time under these linear assumptions. However, for $\theta \approx 0$ we need a more sophisticated model.

A.1.1 Swing Up Control

One way to get the cart to swing the pendulum up to the linear-range is to find a homoclinic orbit (a trajectory that passes through an unstable fixed point). I.e. we must find a controller that drives the pendulum to the unstable equilibrium. This can be done using energy shaping, and in the context of the inverted pendulum, this constitutes applying force to maximise the potential energy and minimise kinetic. Once in the linear region we then switch to an LQR controller to complete the task.

A.2 Propagation of Quantisation Error

The state space model for the quantisation of the linearised inverted pendulum can be written as:

$$\mathbf{x}_t^{(2D)} = \mathbf{C}\mathbf{x}_t + \mathbf{V}_t \quad \mathbf{V}_t \sim WN(0, \sigma_v^2 \mathbf{I}) \quad (\text{A.5})$$

$$\mathbf{x}_t = \mathbf{A}\mathbf{x}_{t-1} + \mathbf{B}\mathbf{u}_t \quad (\text{A.6})$$

Where A and B are the linearised system dynamics (valid for small time steps), and C is the linear transformation to a 2D state space, with quantisation noise \mathbf{V} .

*** This derivation here is for 1D and also follows the derivation from 4F7, which does not take into account actions ($B=0$), and has additional noise on the second equation.

However, because actions are deterministic, it should have no effect on the propagated MSE. Could fairly easily do a full derivation for this, but is it that necessary? Also σ_v^2 is not true, as this is only a 1D error. Finally, possible spanner... Limiting the memory of the system to only t-k could place limits on the accuracy? Should do full derivation if time allows. Note, is this even the correct way to go about this? Can we assume a neural network will act as a kalman filter?***

Assuming the quantisation bin size, δx , is small, the quantisation noise can be modelled as uniform random variable with a mean squared error:

$$\sigma_v^2 = \mathbb{E}[V_n^2] = \int_{-\delta x/2}^{\delta x/2} u^2 \cdot \frac{1}{\delta x} du = \frac{\delta x^2}{12} \quad (\text{A.7})$$

Kalman filtering can be used to find an optimal estimate $\hat{X}_n = K[X_n|Y_{n-k:n}]$ using the algorithm (derived in [3]).

Algorithm 3 Kalman Filtering

1: **Given:** $\hat{X}_n = K[X_n|Y_{n-k:n}]$ and $\sigma^2 = \mathbb{E}[(X_n - \hat{X}_n)^2]$

Prediction:

2: $\bar{X}_{n+1} = K[X_n|Y_{n-k:n}] = f_n \hat{X}_n$

3: $\bar{\sigma}_{n+1}^2 = \mathbb{E}[(X_{n+1} - \bar{X}_{n+1})^2] = f_n \sigma_n^2 + \sigma_w^2 \quad \triangleright \sigma_w^2 = 0$

Update:

4: $I_{n+1} = Y_{n+1} - g_{n+1} \bar{X}_{n+1}$

5: $\hat{X}_{n+1} = \bar{X}_{n+1} + \frac{g_{n+1} \bar{\sigma}_{n+1}^2}{g_{n+1}^2 \bar{\sigma}_{n+1}^2 + \sigma_v^2} I_{n+1} \quad \triangleright \sigma_v^2 = \frac{\delta x^2}{12}$

6: $\sigma_{n+1}^2 = \bar{\sigma}_{n+1}^2 \left(1 - \frac{g_{n+1}^2 \bar{\sigma}_{n+1}^2}{g_{n+1}^2 \bar{\sigma}_{n+1}^2 + \sigma_v^2} \right)$

Thus we can find the optimal mean squared error of the next value from algorithm 3 line 6:

$$\sigma_{n+1}^2 = f_n \sigma_n^2 \left(1 - \frac{f_n \sigma_n^2}{f_n \sigma_n^2 + \sigma_v^2} \right) \quad (\text{A.8})$$