

# On Monte Carlo Tree Search With Multiple Objectives



Michael Painter  
Pembroke College  
University of Oxford

A thesis submitted for the degree of  
*Doctor of Philosophy*

Trinity 2024

TODO: list

1. THTS section, as it defines notation. Put macros in text/abbreviations.tex
2. Copy DENTS and BTS into thesis and copy into common notation
3. Define the toy problems, D-Chain and the one with the entropy trap

# Acknowledgements

TODO: acknowledgements here



# Abstract

TODO: abstract here



# Contents

<b>List of Figures</b>	<b>ix</b>
<b>List of Tables</b>	<b>xi</b>
<b>List of Notation</b>	<b>xiii</b>
<b>List of Abbreviations</b>	<b>xvii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Overview . . . . .	1
1.2 Contributions . . . . .	2
1.3 Structure of Thesis . . . . .	4
1.4 Publications . . . . .	4
<b>2 Background</b>	<b>5</b>
2.1 Markov Decision Processes and Reinforcement Learning . . . . .	6
2.1.1 Maximum Entropy Reinforcement Learning . . . . .	10
2.1.2 Remaining todos for this chapter after first draft . . . . .	11
2.2 Trial-Based Heuristic Tree Search and Monte-Carlo Tree Search . .	12
2.2.1 Trial Based Heuristic Tree Search . . . . .	13
2.2.2 Upper Confidence Bounds Applied to Trees (UCT) . . . . .	17
2.2.3 Monte-Carlo Tree Search . . . . .	18
2.2.4 Maximum Entropy Tree Search . . . . .	20
2.3 Multi-Objective Reinforcement Learning . . . . .	24
2.3.1 Convex Hull Value Iteration . . . . .	28
2.4 Multi-Objective Monte Carlo Tree Search . . . . .	28
2.5 Sampling Random Variables . . . . .	29
<b>3 Literature Review</b>	<b>31</b>
3.1 Multi-Armed Bandits . . . . .	31
3.2 Reinforcement Learning . . . . .	32
3.3 Trial-Based Heuristic Tree Search and Monte-Carlo Tree Search . .	32
3.3.1 Trial Based Heuristic Tree Search . . . . .	32

3.3.2	Monte-Carlo Tree Search . . . . .	32
3.3.3	Maximum Entropy Tree Search . . . . .	32
3.4	Multi-Objective Reinforcement Learning . . . . .	32
3.5	Multi-Objective Monte Carlo Tree Search . . . . .	33
3.6	Sampling Random Variables . . . . .	33
<b>4</b>	<b>Monte Carlo Tree Search With Boltzmann Exploration</b>	<b>35</b>
4.1	Introduction . . . . .	35
4.2	Boltzmann Search . . . . .	36
4.3	Toy Environments . . . . .	36
4.4	Theoretical Results . . . . .	36
4.5	Empirical Results . . . . .	36
4.6	Full Results . . . . .	37
<b>5</b>	<b>Convex Hull Monte Carlo Tree Search</b>	<b>39</b>
5.1	Introduction . . . . .	39
5.2	Contextual Tree Search . . . . .	39
5.3	Contextual Zooming for Trees . . . . .	40
5.4	Convex Hull Monte Carlo Tree Search . . . . .	40
5.5	Results . . . . .	40
<b>6</b>	<b>Simplex Maps for Multi-Objective Monte Carlo Tree Search</b>	<b>41</b>
6.1	Introduction . . . . .	41
6.2	Simplex Maps . . . . .	42
6.3	Simplex Maps in Tree Search . . . . .	42
6.4	Theoretical Results . . . . .	42
6.5	Empirical Results . . . . .	42
<b>7</b>	<b>Conclusion</b>	<b>43</b>
7.1	Summary of Contributions . . . . .	43
7.2	Future Work . . . . .	43
<b>Appendices</b>		
<b>A</b>	<b>List Of Appendices To Consider</b>	<b>47</b>
<b>Bibliography</b>		<b>49</b>



## List of Figures



## List of Tables



# List of Notation

$\mathbb{1}$  . . . . . The indicator function, where  $\mathbb{1}(A) = 1$  when  $A$  is true, and  $\mathbb{1}(A) = 0$  when  $A$  is false.

## Markov Decision Processes and Reinforcement Learning (S)

$\alpha$  . . . . . The temperature parameter. (The coefficient of the entropy term in the maximum entropy (soft) objective).

$\mathcal{A}$  . . . . . Set of actions in an MDP.

$H$  . . . . . The finite-horizon time bound of an MDP.

$\mathcal{H}$  . . . . . Shannon entropy, of a policy.

$J(\pi)$  . . . . . **TODO: objective function**

$J_{\text{sft}}(\pi)$  . . . . . **TODO: soft objective function**

$p$  . . . . . The next-state probability distribution of an MDP.  $p(\cdot|s, a) : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ .

$\pi$  . . . . . A policy, **TODO: state that this is just for this chapter, see sect 2.2 defn for rest of thesis**

$\pi^*$  . . . . . The optimal standard policy, that maximises the objective function  $J(\pi)$

$\pi_{\text{sft}}^*$  . . . . . The optimal soft policy, that maximises the soft objective function  $J_{\text{sft}}(\pi)$

$Q^\pi$  . . . . . The Q-value of a policy  $\pi$ .  $Q^\pi(s, a; t)$  denotes the expected cumulative reward that policy  $\pi$  will achieve, starting from state  $s_t = s$ , starting by taking action  $a_{t+1} = a$ .

$Q^*$  . . . . . The optimal value function  $Q^*(s, a; t)$  denotes the expected maximal value that can be achieved from state  $s_t = s$ , starting with action  $a_{t+1} = a$  by any policy.

$Q_{\text{sft}}^\pi$  . . . . . The Q-value of a policy  $\pi$ .  $Q_{\text{sft}}^\pi(s, a; t)$  denotes the expected cumulative reward that policy  $\pi$  will achieve, starting from state  $s_t = s$ , starting by taking action  $a_{t+1} = a$ , **TODO: plus the entropy of the policy weighted by the temperature alpha. TODO: consider just writing max pi over the soft values here instead?**

$Q_{\text{sft}}^*$ . . . . .	The optimal value function $Q_{\text{sft}}^*(s, a; t)$ denotes the expected maximal value that can be achieved from state $s_t = s$ , starting with action $a_{t+1} = a$ by any policy, <b>TODO: plus the entropy of the (optimal soft) policy weighted by the temperature alpha.</b>
$R$ . . . . .	The reward function of in an MDP: $\mathcal{R}(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ .
$\mathcal{S}$ . . . . .	Set of states in an MDP.
$\tau$ . . . . .	<b>TODO: trajectory</b>
$\tau_{:h}$ . . . . .	<b>TODO: trajectory</b>
$t$ . . . . .	<b>TODO: free variable for current timestep? And generally add a list of free variables?</b>
$V^\pi$ . . . . .	The value of a policy $\pi$ . $V^\pi(s; t)$ denotes the expected cumulative reward that policy $\pi$ will achieve, starting from state $s_t = s$ .
$V^*$ . . . . .	The optimal value function $V^*(s; t)$ denotes the expected maximal value that can be achieved from state $s_t = s$ by any policy.
$V_{\text{sft}}^\pi$ . . . . .	The soft value of a policy $\pi$ . $V_{\text{sft}}^\pi(s; t)$ denotes the expected cumulative reward that policy $\pi$ will achieve, starting from state $s_t = s$ , <b>TODO: plus the entropy of the policy weighted by the temperature alpha.</b>
$V_{\text{sft}}^*$ . . . . .	The optimal soft value function $V^*(s; t)$ denotes the expected maximal value that can be achieved from state $s_t = s$ by any policy, <b>TODO: plus the entropy of the (optimal soft) policy weighted by the temperature alpha.</b> <b>TODO: consider just writing max pi over the soft values here instead?</b>
$\mathcal{M}$ . . . . .	<b>TODO: define here and move into correct place in the list</b>

## Trial Based Heuristic Tree Search (Section 2.2)

$\mathcal{B}_Q$ . . . . .	<b>TODO: define backup</b>
$\mathcal{B}_V$ . . . . .	<b>TODO: define backup</b>
$N(s)$ . . . . .	The number of visits at the decision node corresponding to state $s$ .
<code>mcts_mode</code> . . . . .	<b>TODO: definition of <code>mcts_mode</code></b>
$N(s)$ . . . . .	The number of visits at the decision node corresponding to state $s$ .
$N(s, a)$ . . . . .	The number of visits at the chance node corresponding to state-action pair $(s, a)$ .

$n$ . . . . .	Number of trials run.
<b>node</b> . . . . .	A mapping from states and state-action pairs to their corresponding decision and chance nodes respectively.
$\text{node}(s).V$ . . .	TODO: write this
$\text{node}(s, a).Q$ . .	TODO: write this
$\pi$ . . . . .	A search policy, TODO: define, this is a parameter of thts++. TODO: handle $\pi^k$
$T$ . . . . .	Computation time limit.
$\mathcal{T}$ . . . . .	A THTS search tree. TODO: With: $\mathcal{T} \subseteq \mathcal{S} \cup \mathcal{S} \times \mathcal{A}$ . TODO: handle $\mathcal{T}^k$
$\mathcal{T}^k$ . . . . .	TODO: search tree on after k trials
$V_{\text{init}}$ . . . . .	The initialisation function used in THTS++, used to initialise the value of a new decision node.
$\bar{V}_{\text{UCT}}$ . . . . .	TODO: define here and move into correct place in the list
$\bar{Q}_{\text{UCT}}$ . . . . .	TODO: define here and move into correct place in the list
$\pi_{\text{UCT}}$ . . . . .	TODO: define here and move into correct place in the list
$b_{\text{UCT}}$ . . . . .	TODO: define here and move into correct place in the list
$\bar{V}_{\text{MCTS}}$ . . . . .	TODO: define here and move into correct place in the list
$\bar{Q}_{\text{MCTS}}$ . . . . .	TODO: define here and move into correct place in the list
$\pi_{\text{MCTS}}$ . . . . .	TODO: define here and move into correct place in the list
$b_{\text{MCTS}}$ . . . . .	TODO: define here and move into correct place in the list
$\hat{V}_{\text{MENTS}}$ . . . . .	TODO: define here and move into correct place in the list
$\hat{Q}_{\text{MENTS}}$ . . . . .	TODO: define here and move into correct place in the list
$\pi_{\text{MENTS}}$ . . . . .	TODO: define here and move into correct place in the list
$\alpha_{\text{MENTS}}$ . . . . .	TODO: define here and move into correct place in the list
$\lambda_s$ . . . . .	TODO: define here and move into correct place in the list
$\epsilon$ . . . . .	TODO: define here and move into correct place in the list

## Multi-Objective Reinforcement Learning (Section 2.3)

$\mathcal{M}$ . . . . .	TODO: define here and move into correct place in the list
$\tau$ . . . . .	TODO: define here and move into correct place in the list
$\tau_{:h}$ . . . . .	TODO: define here and move into correct place in the list

$\mathbf{V}^\pi$ . . . . .	TODO: define here and move into correct place in the list
$\mathbf{Q}^\pi$ . . . . .	TODO: define here and move into correct place in the list
$\Delta^d$ . . . . .	TODO: define here and move into correct place in the list
$u$ . . . . .	TODO: define here and move into correct place in the list
$\mathbf{w}$ . . . . .	TODO: define here and move into correct place in the list
$d$ . . . . .	TODO: define here and move into correct place in the list
$\Pi$ . . . . .	TODO: define here and move into correct place in the list
$U(\Pi; u)$ . . . . .	TODO: define here and move into correct place in the list
$u_{\text{lin}}$ . . . . .	TODO: define here and move into correct place in the list
$CH(\Pi)$ . . . . .	TODO: define here and move into correct place in the list
$CS(\Pi; u)$ . . . . .	TODO: define here and move into correct place in the list
$CCS(\Pi)$ . . . . .	TODO: define here and move into correct place in the list
. . . . .	TODO: define here and move into correct place in the list
. . . . .	TODO: define here and move into correct place in the list
. . . . .	TODO: define here and move into correct place in the list



## List of Abbreviations

### Markov Decision Processes and Reinforcement Learning (S)

**MDP** . . . . . Markov Decision Process.

### Trial Based Heuristic Tree Search (Section 2.2)

**MCTS** . . . . . Monte Carlo Tree Search.

**MENTS** . . . . . Maximum ENtropy Tree Search.

**THTS** . . . . . Trial-based Heuristic Tree Search.

**THTS++** . . . . . **TODO: thts++**

**UCT** . . . . . Upper Confidence Bound applied to Trees.

**node** . . . . . **TODO: define here and move into correct place in the list**

**MOMDP** . . . . Multi-Objective Markov Decision Process



# 1

## Introduction

### Contents

---

1.1	Overview . . . . .	1
1.2	Contributions . . . . .	2
1.3	Structure of Thesis . . . . .	4
1.4	Publications . . . . .	4

---

TODO: chapter structure (i.e. in the introduction section I give some background in the field(s), cover the main contributions of this thesis, etc, etc).

## 1.1 Overview

TODO: list

- Give some context around MCTS (and talk about exploration and exploitation), and why we might use it
  - Larger scale than tabular methods
  - Can do probability and theory stuff (and some explainability, by looking at stats in the tree the agent used)
  - Can use tree search with neural networks to get some of the above (and use for neural network training as in alpha zero)

- Argument from DENTS paper for exploration  $>$  exploitation (in context of planning in a simulator)
- Give high level overview of Multi-Objective RL, and why it can be useful
- Give an idea of how my work fits into MCTS and MORL as a whole
- Discuss research questions/issues with current literature (i.e. introduce some of the ideas from contributions section below)

## 1.2 Contributions

TODO: Inline acronyms used, or make sure that they're defined before hand

Throughout this thesis, we will consider the following questions related to Monte Carlo Tree Search and Multi-Objective Reinforcement Learning:

**Q1 - Exploration:** When planning in a simulator with limited time, how can MCTS algorithms best explore to make good decisions?

**Q1.1 - Entropy:** Entropy is often used as an exploration objective in RL, but can it be used soundly in MCTS?

**Q1.2 - Multi-Objective Exploration:** How can Multi-Objective MCTS methods explore to find optimal actions for different objectives?

**Q2 - Scalability:** How can the scalability of (multi-objective) MCTS methods be improved?

**Q2.1 - Complexity:** MCTS algorithms typically run in  $O(nAH)$ , but are there algorithms that can improve upon this?

**Q2.2 - Multi-Objective Scalability:** With respect to the size of environments, how scalable are Multi-Objective MCTS methods?

**Q2.3 - Curse of Dimensionality:** With respect to the number of objectives, to what extent do Multi-Objective MCTS methods suffer from the curse of dimensionality?

**Q3 - Evaluation:** How can we best evaluate a search tree produced by a Monte Carlo Tree Search algorithm?

**Q3.1 - Tree Policies:** Does it suffice to extract a policy from a single search tree for evaluation? **TODO:** going to have to run some extra experiments for that, but I probably should do that for completeness anyway

**Q3.2 - Multi-Objective Evaluation:** Can we apply methods from the MORL literature to theoretically and empirically evaluate Multi-Objective MCTS?

**TODO:** some words about how below is the contributions we're making in this thesis and expand these bullets a bit more

- Max Entropy can be misaligned with reward maximisation (**Q1.1 - Entropy**)
- Boltzmann Search Policies - BTS and DENTS (**Q1.1 - Entropy**, and with extra results **Q3.1 - Tree Policies**)
- Use the alias method to make faster algorithms (**Q2.1 - Complexity**)
- Simple regret (**Q1 - Exploration**)
- Use of contexts in THTS to make consistent decisions in each trial (**Q1.2 - Multi-Objective Exploration**)
- Contextual regret introduced in CHMCTS (**Q2.2 - Multi-Objective Scalability**, **Q3.2 - Multi-Objective Evaluation**)
- Contextual Zooming and CHMCTS (designed for **Q1.2 - Multi-Objective Exploration**, runtimes cover **Q2.3 - Curse of Dimensionality**, results **Q3.2 - Multi-Objective Evaluation**)
- Simplex maps (**Q1.2 - Multi-Objective Exploration**, **Q2.2 - Multi-Objective Scalability**, **Q2.3 - Curse of Dimensionality**)
- Contextual Simple Regret (**Q3.2 - Multi-Objective Evaluation**)

TODO: Would like to do the comparing different types of eval, even if not listing it as a research question (compare giving it X seconds per decision and evaluating that policy (SLOW), and comparing policy extracted from the tree)

TODO: can make an argument that the best bound achieved by theory is given by letting temperature go to max. Which is consistent with the exploring bandits results

## 1.3 Structure of Thesis

TODO: a paragraph with a couple lines to a paragraph about each chapter. This is the high level overview/intro to the thesis paragraph. I.e. this section is “this is the story of my thesis in a page or two”

## 1.4 Publications

TODO: update final publication when submit

The work covered in this thesis also appears in the following publications:

- Painter, M; Lacerda, B; and Hawes, N. “Convex Hull Monte-Carlo Tree-Search.” In *Proceedings of the international conference on automated planning and scheduling. Vol. 30. 2020*, ICAPS, 2020, (see Appendix ??).
- Painter, M; Baïoumy, M; Hawes, N; and Lacerda, B. “Monte Carlo Tree Search With Boltzmann Exploration.” In *Advances in Neural Information Processing Systems, 36, 2023*, NeurIPS, 2023, (see Appendix ??).
- Painter, M; Hawes, N; and Lacerda, B. “Simplex Maps for Multi-Objective Monte Carlo Tree Search.” In *TODO, Under Review at conf\_name*, (see Appendix ??).

# 2

## Background

### Contents

---

<b>2.1</b>	<b>Markov Decision Processes and Reinforcement Learning</b>	<b>6</b>
2.1.1	Maximum Entropy Reinforcement Learning . . . . .	10
2.1.2	Remaining todos for this chapter after first draft . . . . .	11
<b>2.2</b>	<b>Trial-Based Heuristic Tree Search and Monte-Carlo Tree Search</b>	<b>12</b>
2.2.1	Trial Based Heuristic Tree Search . . . . .	13
2.2.2	Upper Confidence Bounds Applied to Trees (UCT) . . . . .	17
2.2.3	Monte-Carlo Tree Search . . . . .	18
2.2.4	Maximum Entropy Tree Search . . . . .	20
<b>2.3</b>	<b>Multi-Objective Reinforcement Learning</b>	<b>24</b>
2.3.1	Convex Hull Value Iteration . . . . .	28
<b>2.4</b>	<b>Multi-Objective Monte Carlo Tree Search</b>	<b>28</b>
<b>2.5</b>	<b>Sampling Random Variables</b>	<b>29</b>

---

TODO: Introduce that going to introduce notation and give the building blocks this thesis builds off

TODO: Nicks comment to revisit: In Sections 2 and 3 you present RL before tree search. This was a surprise to me. Wouldn't it be better to start with single objective MDPs, then introduce methods of solving them (planning, bandits, rl) and the different assumptions they make? Then move to MO versions. I'd also be open to you presenting MCTS then generalising to MCTS rather than the other way around, but I think the way you present it is probably the most efficient.

TODO: Add a regret or multi-armed bandit section?

## 2.1 Markov Decision Processes and Reinforcement Learning

TODO: list

- Typical agent interacting with environment diagram
- Agent planning with simulator
- MDPs definition
- Policies
- Value functions (single (and multi-objective?))
- Basic results and definitions we use (tabular planning algorithms)
- Talk about entropy and some of that work (probably a subsection)

TODO: Split this section into MDPs definition and move RL below monte carlo tree search? Probably if do the multi-armed bandit stuff.

In this section Markov Decision Processes are introduced, as well as some fundamental concepts in Reinforcement Learning such as *policies* and *value functions* which will be useful in the following. Afterwards, the *maximum entropy objective* is discussed, as it will be relevant to some of the work in this thesis and other closely related work.

In this thesis we will only consider discrete and finite-horizon Markov Decision Processes, TODO: but we will point to works that make adaptations for non-finite horizons and continuous and partially observable environments, and briefly describe why the ideas could be used in parallel with the concepts in this thesis.

**Definition 2.1.1.** A Markov Decision Process (MDP) is a tuple  $\mathcal{M} = (\mathcal{S}, s_0, \mathcal{A}, R, p, H)$ , where  $\mathcal{S}$  is a set of states,  $s_0 \in \mathcal{S}$  is an initial state,  $\mathcal{A}$  is a set of actions,  $R(s, a)$  is a reward function  $\mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ ,  $p(s'|s, a)$  is a next state transition distribution  $\mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  and  $H \in \mathbb{N}$  is a finite-horizon time bound.



Where it is more convenient, we will use the following notation to refer to the set of successor states when discussing MDPs:

**Definition 2.1.2.** *The set of successor states  $\text{Succ}(s, a)$  of a state-action pair  $(s, a)$  is defined as  $\text{Succ}(s, a) = \{s' | p(s'|s, a) > 0\}$ . Additionally, let  $s' \sim \text{Succ}(s, a)$  be a shorthand for  $s' \sim p(\cdot|s, a)$ .*

TODO: decide if should keep  $p(\cdot|s, a)$ , and just swap out the succ command as shorthand for that?

Formally, we will consider a policy to be a mapping from a state in  $\mathcal{S}$  to a distribution over actions  $\mathcal{A}$ . TODO: For this thesis we will consider policies to be stochastic, so when a policy is used to generate actions to follow, it is done by sampling from the distribution.

TODO: So im using  $\pi$  for the definition here. But for rest of thesis it will be used as explicitly the "search policy"

In some cases it will be necessary to talk about deterministic policies. TODO: When ambiguous we will state when a policy is deterministic. A deterministic policy is written as a one hot stochastic distribution

**Definition 2.1.3.** *A (stochastic) policy  $\pi : \mathcal{S} \rightarrow \mathcal{A} \rightarrow [0, 1]$  is a mapping from states to distributions over actions, where for all  $s \in \mathcal{S}$  we have  $\sum_{a \in \mathcal{A}} \pi(a|s) = 1$ . We will use the following notations:  $\pi(\cdot|s)$  is used to denote the entire distribution over the set of actions,  $\pi(a|s)$  denotes the probability that action  $a$  is sampled by the policy at state  $s$ .*

Additionally, for deterministic policies we will use the shorthand notation  $\pi(s)$  to denote the action taken by the policy. Formally, we will write  $\pi(s) = a'$  where  $a' \in \mathcal{A}$  is such that  $\pi(a'|s) = 1$ .

TODO: want to add a statement saying something like  $\pi(s) = a'$  is shorthand for  $\pi(a|s) = \mathbb{1}[a = a']$ .

TODO: Actually want to say something like: a deterministic policy is special case of the stochastic, with a delta distribution, such as  $\pi(a|s) = \mathbb{1}[a = a']$ , and in

such cases we will use  $\pi(s) = a'$  as a short hand to state  $\pi(a|s) = \mathbb{1}[a = a']$  and  $\pi(s)$  as a shorthand for  $a'$ .

TODO: Explicitly note that we are using  $\pi$  as a free variable (is free variable correct word here?) here, and that in section - todo ref - we will use  $\pi$  to refer to a specific type of policy (search policies) in the rest of the thesis after this section.

TODO: clean up policy definition with shorthand stuff

TODO: actually state the the values of a policy need to add to one

In reinforcement learning it is common to refer to a sequence of states, actions and rewards as a *trajectory*. In the tree search literature it is more common to refer to this as a *trial* TODO: cite thts?, so we will use these terms interchangeably.

**Definition 2.1.4.** A trajectory, is a sequence of state, action and rewards, that is induced by a policy  $\pi$  and MDP  $\mathcal{M}$  pair. Let the trials/trajectory be  $\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_{H-1}, a_{H-1}, r_{H-1}, s_H)$ , where  $a_t \sim \pi(\cdot|s_t)$ ,  $r_t = R(s_t, a_t)$  and  $s_{t+1} \sim \text{Succ}(s_t, a_t)$ . Notationally, we will write  $\tau \sim \pi$  to denote a sampled trial/trajectory with respect to a policy, where the MDP is implicit.

Sometimes it will be necessary to reason about trajectories with a horizon  $h < H$ , which will be denoted  $\tau_h = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_{h-1}, a_{h-1}, r_{h-1}, s_h)$ .

Next the value and Q-value of a policy is defined, which is the expected cumulative reward that a policy will obtain:

**Definition 2.1.5.** The value of a policy  $\pi$  from state  $s$  at time  $t$  is:

$$V^\pi(s; t) = \mathbb{E}_{\tau \sim \pi} \left[ \sum_{i=t}^{H-1} r_i \middle| s_t = s \right]. \quad (2.1)$$

The Q-value of a policy  $\pi$ , from state  $s$ , with action  $a$ , at time  $t$  is:

$$Q^\pi(s, a; t) = R(s, a) + \mathbb{E}_{s' \sim \text{Succ}(s, a)} [V^\pi(s'; t + 1)]. \quad (2.2)$$

From the definition of the values functions the optimal value functions can be defined:

**Definition 2.1.6.** *The Optimal (Q-)Value of a state(-action pair) is defined as:*

$$V^*(s; t) = \max_{\pi} V^{\pi}(s; t) \quad (2.3)$$

$$Q^*(s, a; t) = \max_{\pi} Q^{\pi}(s, a; t). \quad (2.4)$$

In reinforcement learning, the objective is to find a policy with maximal value:

**Definition 2.1.7.** *The (standard) reinforcement learning objective function  $J(\pi)$  is defined as:*

$$J(\pi) = V^{\pi}(s_0; 0). \quad (2.5)$$

*The objective of (standard) reinforcement learning can then be stated as finding  $\max_{\pi} J(\pi)$ .*

It can be shown [TODO: refs](#) that the optimal (Q-)value functions satisfy the *Bellman equations*:

$$V^*(s; t) = \max_{a \in \mathcal{A}} Q^*(s, a; t), \quad (2.6)$$

$$Q^*(s, a; t) = R(s, a) + \mathbb{E}_{s' \sim \text{Succ}(s, a)} [V^*(s'; t + 1)]. \quad (2.7)$$

In tabular reinforcement learning, a table of values for  $V(s; t)$  [TODO: havent actually defined V without any superscript](#) is kept for each  $s, t$ . Given any initial value function  $V^{(0)}$  let the *Bellman backup* operations be:

$$V^{k+1}(s; t) = \max_{a \in \mathcal{A}} Q^{k+1}(s, a; t), \quad (2.8)$$

$$Q^{k+1}(s, a; t) = \mathbb{E}_{s' \sim \text{Succ}(s, a)} [R(s, a) + V^k(s'; t + 1)]. \quad (2.9)$$

Using this *dynamic programming* approach is known as *value iteration*. It can be shown that the Bellman backups are contraction operators [TODO: add cite](#), which can be used to show that  $V^k \rightarrow V^*$  as  $k \rightarrow \infty$ . In the discrete [TODO: what the actual conditions are](#) case we are considering, there will always be some  $N < \infty$  such that  $V^N = V^*$ .

[TODO: add optimal policy from Q values](#)

### 2.1.1 Maximum Entropy Reinforcement Learning

In *Maximum Entropy Reinforcement Learning*, the objective function is altered to include the addition of an entropy term. Let  $\mathcal{H}$  denote the (Shannon) entropy function TODO: cite:

$$\mathcal{H}(\pi(\cdot|s)) = \mathbb{E}_{a \sim \pi(\cdot|s)}[-\log \pi(a|s)]. \quad (2.10)$$

Note that there are other forms of entropy, such as relative and Tsallis entropy, which can be used in place of Shannon entropy TODO: cite. For the work considered in this thesis, the other forms of entropy can be used by replacing the definition of  $\mathcal{H}$  by the relevant definition.

In the maximum entropy objective, the relative weighting of entropy terms is included using a coefficient  $\alpha$ , which is called the *temperature*. In the maximum entropy objective, analogues of the value functions can be defined, which are typically referred to as *soft (Q-)values*, and similarly the maximum entropy objective is often referred to as the *soft objective*.

**Definition 2.1.8.** *The soft value of a policy  $\pi$  from state  $s$  at time  $t$  is:*

$$V_{\text{sft}}^\pi(s; t) = \mathbb{E}_{\tau \sim \pi} \left[ \sum_{i=t}^{H-1} r_t + \alpha \mathcal{H}(\pi(\cdot|s_i)) \middle| s_t = s \right]. \quad (2.11)$$

*The soft Q-value of a policy  $\pi$ , from state  $s$ , with action  $a$ , at time  $t$  is:*

$$Q_{\text{sft}}^\pi(s, a; t) = R(s, a) + \mathbb{E}_{s' \sim p(\cdot|s, a)}[V_{\text{sft}}^\pi(s'; t + 1)]. \quad (2.12)$$

Similarly, optimal soft (Q-)values can be defined:

**Definition 2.1.9.** *The Optimal soft (Q-)Value of a state(-action pair) is defined as:*

$$V_{\text{sft}}^*(s; t) = \max_{\pi} V_{\text{sft}}^\pi(s; t) \quad (2.13)$$

$$Q_{\text{sft}}^*(s, a; t) = \max_{\pi} Q_{\text{sft}}^\pi(s, a; t). \quad (2.14)$$

Equations similar to the Bellman equations, aptly named the *Soft Bellman equations* can be defined, which differ to equations [TODO: ref](#) by the replacement of the max operation with the *softmax* operation (which is why the maximum entropy analogues are referred to as the *soft* versions of their standard reinforcement learning counterparts).

In maximum entropy reinforcement learning, the objective is to find a policy with maximal soft value:

**Definition 2.1.10.** *The maximum entropy (or soft) reinforcement learning objective function  $J_{\text{sft}}(\pi)$  is defined as:*

$$J_{\text{sft}}(\pi) = V_{\text{sft}}^{\pi}(s_0; 0). \quad (2.15)$$

*The objective of maximum entropy (or soft) reinforcement learning can then be stated as finding  $\max_{\pi} J_{\text{sft}}(\pi)$ .*

Similarly to standard reinforcement learning, it can be shown [TODO: refs](#) that the optimal soft (Q-)value functions satisfy the *soft Bellman equations*:

$$V_{\text{sft}}^*(s; t) = \alpha \log \sum_{a \in \mathcal{A}} \exp(Q_{\text{sft}}^*(s, a; t) / \alpha), \quad (2.16)$$

$$Q_{\text{sft}}^*(s, a; t) = R(s, a) + \mathbb{E}_{s' \sim \text{Succ}(s, a)}[V_{\text{sft}}^*(s'; t + 1)]. \quad (2.17)$$

Again, similarly to standard reinforcement learning, we can define *soft Bellman backups* that admit an analogous algorithm to value iteration:

$$V_{\text{sft}}^{k+1}(s; t) = \alpha \log \sum_{a \in \mathcal{A}} \exp(Q_{\text{sft}}^{k+1}(s, a; t) / \alpha), \quad (2.18)$$

$$Q_{\text{sft}}^{k+1}(s, a; t) = R(s, a) + \mathbb{E}_{s' \sim \text{Succ}(s, a)}[V_{\text{sft}}^k(s'; t + 1)]. \quad (2.19)$$

Finally, given the optimal soft value and soft Q-value functions, the optimal soft policy is known [TODO: cite](#):

$$\pi_{\text{sft}}^*(a|s; t) = \exp((Q_{\text{sft}}^*(s, a; t) - V_{\text{sft}}^*(s; t)) / \alpha). \quad (2.20)$$

### 2.1.2 Remaining todos for this chapter after first draft

[TODO:](#) Add a comment similar to DENTS paper where we will drop the  $t$  in notation. Make it quite bold somehow

## 2.2 Trial-Based Heuristic Tree Search and Monte-Carlo Tree Search

TODO: list

- Give high level overview of MCTS (why use it etc)
- Outline that I'll present this as here is THTS, and then here's the THTS routines for MCTS

TODO: double check the intro to 2.2, as wrote this a while ago

TODO: try to make sure specific about using thts vs mcts

In this section we introduce **THTS++** [3], which is an open-source, parallelised extension of the Trial-based Heuristic Tree Search schema [2] (THTS). This schema is a generalisation of Monte Carlo Tree Search (MCTS), as presented in Section ???. In **THTS++** trees consist of *decision nodes* and *chance nodes*. Decision nodes output actions that can be taken by the agent, and chance nodes output *outcomes* that may be random and may depend on the action taken. As such, each decision node has an associated *state* and each chance node has an associated *state-action pair*. In this work, we are considering fully-observable environments, but **THTS++** can be generalised to consider *partially-observable* environments. We give **THTS++** implementations of the standard Upper Confidence Bound applied to Trees (UCT) algorithm and Maximum ENtropy Tree Search (MENTS) in Sections ?? and ?? respectively.

In MCTS we run trials, either for some fixed number of trials, or some timelimit, where each trial is split into four stages: (1) selection, which samples states and actions for the trial, corresponding to a path down the tree; (2) expansion, which creates any new nodes in the tree; (3) initialisation, which initialises values at any new leaf nodes in the tree; (4) backup, which updates values at all nodes visited on the trial.

TODO: add MCTS figure here?

```

1 def foo(bar):
2     print("helloworld!")

```

## 2.2.1 Trial Based Heuristic Tree Search

TODO: list

- Copy DENTS MCTS section presentation, make a notation  $\text{node}(s_t)$  for the node at state  $s_t$
- Present thts++
- Indicate what parts are new versus the original paper (context function, optionally running `mcts_mode` and multi-threading)
- Small comment about multi-threading and two-phase locking used to avoid deadlock
- TODO: probably not necessary to say - but thought of nice/concise way of explaining it (a node can lock children, not parent, if need info from parent, then it has to put a thread safe copy in the context)
- Define terms precisely and consistently, for example `mcts_mode` (say that notation and terminology varies widely in literature, e.g. does uct run in mcts mode or not?)
- Mention that  $V_{\text{init}}$  can be implemented as  $V_\theta$  to be used with deep RL methods
- TODO: Find the best place to talk about deep RL? Maybe in the RL section?

TODO: this section really really needs diagrams

TODO: add psuedo code

In this section we will present THTS++ schema, which is TODO: adaptation? of the Trial-Based Heuristic Tree Search (THTS) schema TODO: cite. After we have presented THTS++, we will use the schema to define tree search algorithms that are

relevant in this thesis, namely Upper Confidence Bound Applied to Trees (UCT) TODO: cite in subsection TODO: ref and Maximum ENtropy Tree Search (MENTS) in subsection TODO: Ref. Finally we will briefly point out the differences between THTS++ and the original THTS schema in subsection TODO: ref.

TODO: this is already a subsection, so update above

In THTS++ a search tree  $\mathcal{T}$  is built using Monte Carlo trials. Each trial is split into two phases: the *selection phase* where a trajectory is sampled using a *search policy*; and the *backup phase* where value estimates stored in the tree structure are updated. In THTS++ the selection phase encompasses the expansion and initialisation phases TODO: of the common presentation of MCTS, where new nodes are added to the tree and the values of any new leaf node is initialised.

TODO: be more precise about trajectory vs trial, and update for sect 2.1.

To simplify notation in the presentation of THTS++ we will assume that states and state-action pairs have a one to one correspondance with nodes in the search tree  $\mathcal{T}$ . This assumption is purely to simplify notation for a clean presentation, and any results discussed in this thesis generalise to when this assumption does not hold. Given this assumption, we can state that the search tree is a subset of the state and state-action spaces, that is  $\mathcal{T} \subseteq \mathcal{S} \cup \mathcal{S} \times \mathcal{A}$ . TODO: rephrase last sentence for defn

**Definition 2.2.1.** A search tree  $\mathcal{T}$  is a subset of the state and state-action spaces, that is  $\mathcal{T} \subseteq \mathcal{S} \cup \mathcal{S} \times \mathcal{A}$ , where for each  $s \in \mathcal{T}$ , there exists some trajectory  $\tau_{:h}$  such that  $s_h = s$ , each  $s' \in \tau_{:h}$  is also in the tree  $s' \in \mathcal{T}$  and each  $s', a' \in \tau_{:h}$  is also in the tree  $(s', a') \in \mathcal{T}$ .

TODO: clean above defn up

TODO: probably want to explicitly define what it means for s or s,a to be in a trajectory

TODO: words about decision and chance nodeS?

**Definition 2.2.2.** A decision node refers to any state that is in the search tree:  $s \in \mathcal{T}$ . A chance node refers to any state-action pair that is in the search tree:  $(s, a) \in \mathcal{T}$ . And a node is used to refer to any decision or chance node in the tree.



When it is not clear from context if an  $s$  or  $(s, a)$  refers to a state(-action pair), the notation  $\text{node}(s)$  and  $\text{node}(s, a)$  will be used. *TODO: fix node notation here*

Additionally, each decision and chance node will generally store value estimates that are algorithm dependent. To specify this we will use  $\text{node}(s).V$  to denote the set of values stored at node  $\text{node}(s)$ , and  $\text{node}(s, a).Q$  for the set of value stored at node  $\text{node}(s, a)$ . *TODO: make this a defn?*

*TODO: define  $N(s)$  and  $N(s, a)$*

The initial search tree consists of a single root node that corresponds to the initial state of the MDP:  $\mathcal{T}^0 = \{s_0\}$ . And let  $\mathcal{T}^k$  denote the search tree of THTS++ after  $k$  trials have been run.

To specify an algorithm in the THTS++ schema, the following need to be provided:

**Search policy:** A distribution  $\pi^k$  for the  $(k + 1)$ th trial, which can use values in the current search tree  $\mathcal{T}^k$ ;

**Heuristic function:** A function  $V_{\text{init}}$  used as a heuristic to initialise values for new decision nodes added to the tree;

**Backup function:** Two functions  $\mathcal{B}_V$  and  $\mathcal{B}_Q$  which updates values; *TODO: Clean this up, and word better pls*

**MCTS mode:** A boolean `mcts_mode` specifying if THTS++ should operate in MCTS mode.

*TODO: actually define the above things properly somewhere*

The  $k + 1$ th trial of the THTS++ schema operates as follows: *TODO: this is probably better written as psuedocode...*

1. sample a trajectory  $\tau_{:h}$  using the search policy  $\pi^k$ ;
  - If `mcts_mode` is False, then  $h = H$ ;
  - If `mcts_mode` is True, then  $h$  is such that  $s_{h-1} \in \mathcal{T}^k$  and  $s_h \notin \mathcal{T}^k$ , or  $h = H$ .

2. Any new nodes nodes that need to be added from this trajectory are added to the tree,  $\mathcal{T}^{k+1} = \mathcal{T}^k \cup \tau_{:h}$ ;
3. If  $s_h \notin \mathcal{T}^k$  then  $\text{node}(s_h).V$  is initialised using  $V_{\text{init}}$ ;
4. The backup functions are used to update values in the tree:
  - For  $i = h - 1, h - 2, \dots, 1, 0$ :
    - $\text{node}(s_i, a_i).Q \leftarrow \mathcal{B}_Q(\{\text{node}(s').V \mid s' \in \text{node}(s_i, a_i).\text{children}\})$
    - $\text{node}(s_i).V \leftarrow \mathcal{B}_V(\{\text{node}(s_i, a').Q \mid a' \in \text{node}(s_i).\text{children}\})$

TODO: converting the above into some psuedocode. Should probably define children as a property of nodes. Should also just state that  $\text{node}(s_i, a_i).Q$  and so on are just scalar values for now. We can make them vectors when needed later. TODO: also make sure define  $N(s)$  and  $N(s, a)$

```

1 def run_trial(search_policy:  $\pi$ , heuristic_fn:  $V_{\text{init}}$ ):
2      $\tau_{:h} = \text{sample\_trajectory}(\pi)$ 
3     if  $s_h \notin \mathcal{T}$ :
4         initialise_values( $s_h$ ,  $V_{\text{init}}$ )
5     for i in {h-1, h-2, ..., 1, 0}:
6         backup_q( $s_i, a_i$ )
7         backup_v( $s_i$ )
8
9 def sample_trajectory(search_policy:  $\pi$ ):
10     pass
11
12 def initialise_values( $s_h$ ,  $V_{\text{init}}$ ):
13     pass
14
15 def backup_q( $s_i$ ,  $a_i$ ):
16     pass
17
18 def backup_v( $s_i$ ):
19     pass

```

TODO: define  $\text{node}(s_i).\text{children}$  and  $\text{node}(s_i, a_i).\text{children}$ .

TODO: Check how thts deals with using value estimates when children dont exist. There should be some form of using heuristic. Maybe this needs to be added to THTSpp todo list

TODO: make sure neurips paper writing integrated (commented out, and below this comment in the .tex) - read the actual pdf for thesis and neurips papers and compare cover same info?

### 2.2.2 Upper Confidence Bounds Applied to Trees (UCT)

TODO: list

- Define UCT here

TODO: I'm feeling ill writing this section, so just going to word vomit this shit out and make it sound not shit later

TODO: add Quct commands like Qments

Upper Confidence Bounds Applied to Trees (UCT) TODO: cite is a commonly used tree search algorithm, which is based on the Upper Confidence Bounds (UCB) TODO: cite both papers algorithm for Multi-Armed Bandit problems TODO: cite original MAB and a review.

In the literature, UCT and MCTS are often used synonymously, however this leaves some of the specifics of the algorithms used as ambiguous. In this thesis, we will present UCT as it was originally presented in TODO: cite. And in subsection TODO: ref we will specify the variant of UCT which is commonly referred to as MCTS.

UCT can be defined using the THTS schema outlined in section TODO: ref as follows:

Firstly, UCT as originally presented is run with `mcts_modeset` to False. As such, all sampled trajectories are sampled until timestep  $H$ , the finite horizon of the MDP.

At each node  $a$  the sampled averages  $\bar{V}_{UCT}$  or  $\bar{Q}_{UCT}$  for value estimates.

The search policy that UCT follows is:

$$\pi_{UCT}(s) = \arg \max_{a \in \mathcal{A}} Q_{UCT}(s, a) + b_{UCT} \sqrt{\frac{\log(N(s))}{N(s, a)}} \quad (2.21)$$

TODO: add labels for equations

In `TODO: ref above eqn`, when  $N(s, a) = 0$  there is a division by zero, which is taken as  $\inf$ , and ties are broken randomly, which effectively implements the “every arm is initialised by pulling it once” `TODO: actually quote the paper, and cite UCT paper.` `TODO: define the bias param`

After sampling a trajectory  $\tau_{:H} \sim \pi_{\text{UCT}}$  are updated as follows:

$$\bar{Q}_{\text{UCT}}(s_t, a_t) \leftarrow \frac{1}{N(s, a)} \left( (N(s, a) - 1) \bar{Q}_{\text{UCT}}(s_t, a_t) + \sum_{i=t}^{H-1} R(s_i, a_i) \right) \quad (2.22)$$

$$\bar{V}_{\text{UCT}}(s_t) \leftarrow \frac{1}{N(s, a)} \left( (N(s, a) - 1) \bar{V}_{\text{UCT}}(s_t, a_t) + \sum_{i=t}^{H-1} R(s_i, a_i) \right) \quad (2.23)$$

`TODO: add labels for equations`

`TODO: Some note about the V values not actually being used in the algorithm`

`TODO: some comment about it can be implemented as backups (copy equations from THTS), but typically implemented as above. OR, just define backup functions to take the trajectory too`

Because UCT is planning in a finite horizon MDP, the heuristic function will only be called on states that are at the time horizon  $H$ . As such, for UCT we can set  $V_{\text{init}}(s) = 0$ .

`TODO: add polynomial UCT here? and or prioritised UCT from alpha go here?`

`TODO: make sure neurips paper writing integrated in UCT and MCTS section (commented out, and below this comment in the .tex) - read the actual pdf for thesis and neurips papers and compare cover same info?`

## 2.2.3 Monte-Carlo Tree Search

`TODO: list`

- Give overview of MCTS
- Give UCT in terms of THTS schema
- Define terms precisely and consistently in terms of THTS functions, maybe `mcts_mode` should go here
- Define the value initialisation of THTS using a rollout policy for MCTS

- Talk about the things that are ambiguous from literature (e.g. people will just say UCT, which originally presented doesn't run in `mcts_mode`, but often assumed it does)
- Should talk about multi-armed bandits here?

TODO: I'm feeling ill writing this section, so just going to word vomit this shit out and make it sound not shit later

TODO: add Qmcts commands like Qments

In this thesis we will refer to any algorithm that only adds one decision node to the search tree on each trial as an MCTS algorithm. That is any THTS algorithm with `mcts_modeset` to `True` is an MCTS algorithm.

In this section we will present what is commonly referred to as MCTS in the literature, where the heuristic function is either in the form of a *rollout*, using a *rollout policy* TODO: cite papers that do this, including some that just call it UCT, or use a function  $V_\theta$  that aims to approximate the true optimal value function  $V^*$  from Equation TODO: ref TODO: cite papers that do this.

These algorithms follow the same

These algorithms use the similar value functions to UCT,  $\bar{V}_{\text{MCTS}}$  or  $\bar{Q}_{\text{MCTS}}$ .

The search policy corresponds to the UCT search policy, using the new  $\bar{Q}_{\text{MCTS}}$  values:

$$\pi_{\text{MCTS}}(s) = \arg \max_{a \in \mathcal{A}} Q_{\text{MCTS}}(s, a) + b_{\text{MCTS}} \sqrt{\frac{\log(N(s))}{N(s, a)}} \quad (2.24)$$

TODO: add labels for equations

A trajectory  $\tau_{:h}$  is sampled until a new decision node not in the tree is reached, as we are now running with `mcts_modeset` to `True`.

If the algorithm uses a function approximation  $V_\theta$ , then it is used directly for the heuristic function  $V_{\text{init}}$ . If a rollout is used for the heuristic function, then the algorithm needs to define a *rollout policy*  $\pi_{\text{rollout}}$ , which is used to sample a Monte Carlo estimate of the value function  $V^{\pi_{\text{rollout}}}$  as follows. The sampled

trajectory  $\tau_{:h} \sim \pi_{\text{MCTS}}$  is extended with the rollout trajectory  $\tau_{h:H} \sim \pi_{\text{rollout}}$  to give the Monte Carlo estimate of the value at  $s_h$ :

$$V^{\pi_{\text{rollout}}}(s_h) \approx \sum_{i=h}^{H-1} r_i. \quad (2.25)$$

Letting  $\tilde{r} = V_{\text{init}}(s_h)$ , the value estimates (or sample averages) are updated as follows:

$$\bar{Q}_{\text{MCTS}}(s_t, a_t) \leftarrow \frac{1}{N(s, a)} \left( (N(s, a) - 1) \bar{Q}_{\text{MCTS}}(s_t, a_t) + \tilde{r} + \sum_{i=t}^{h-1} r_i \right) \quad (2.26)$$

$$\bar{V}_{\text{MCTS}}(s_t) \leftarrow \frac{1}{N(s, a)} \left( (N(s, a) - 1) \bar{V}_{\text{MCTS}}(s_t, a_t) + \tilde{r} + \sum_{i=t}^{h-1} r_i \right) \quad (2.27)$$

TODO: add labels for equations

## 2.2.4 Maximum Entropy Tree Search

TODO: list

- Define MENTS here

TODO: I'm feeling ill writing this section, so just going to word vomit this shit out and make it sound not shit later

Maximum Entropy Tree Search (MENTS) [TODO: cite](#), in contrast to UCT, focuses on the maximum-entropy objective. In its original presentation `mcts_mode` is set to True, and it uses the soft value estimates  $\hat{V}_{\text{MENTS}}$  and  $\hat{Q}_{\text{MENTS}}$ . The MENTS search policy is

$$\pi_{\text{MENTS}}(a|s) = (1 - \lambda_s) \exp \left( \frac{1}{\alpha_{\text{MENTS}}} \left( \hat{Q}_{\text{MENTS}}(s, a) - \hat{V}_{\text{MENTS}}(s) \right) \right) + \frac{\lambda_s}{|\mathcal{A}|}, \quad (2.28)$$

[TODO: add labels for equations](#) where  $\alpha_{\text{MENTS}}$  is the temperature parameter used for Equation [TODO: ref](#) in MENTS, and  $\lambda_s = \min(1, \epsilon / \log(e + N(s)))$ , with  $\epsilon \in (0, \infty)$  is an exploration parameter.

The value estimates are updated using the soft Bellman backups (TODO: ref) as follows:

$$\hat{Q}_{\text{MENTS}}(s_t, a_t) \leftarrow R(s_t, a_t) + \sum_{s' \in \text{Succ}(s, a)} \left( \frac{N(s')}{N(s_t, a_t)} \hat{V}_{\text{MENTS}}(s') \right), \quad (2.29)$$

$$\hat{V}_{\text{MENTS}}(s_t) \leftarrow \alpha \log \sum_{a \in \mathcal{A}} \exp \left( \frac{1}{\alpha} \hat{Q}_{\text{MENTS}}(s_t, a) \right). \quad (2.30)$$

TODO: add labels for equations

TODO: talk about initialisations, for Vinit its the same as MCTS, think about how to integrate Qinit properly into the thesis using below stuff (commented out)

TODO: make sure neurips paper writing integrated (commented out, and below this comment in the .tex) - read the actual pdf for thesis and neurips papers and compare cover same info?

THINGS WROTE before for THTS section **TODO:** read through and see if anything want to keep, otherwise delete



In THTS++ we run trials for either some fixed number of trials  $n$ , or some time limit  $T$ . Each trial consists of three steps: (1) sample a context, which is used to store variables that are associated with a specific trial, and is passed to the following three functions; (2) selection, which samples states, actions and outcomes for the trial, corresponding to a path down the tree; (3) initialisation, which creates any new nodes in the tree and initialises their values; (4) backup, which updates values at all nodes visited on the trial.

Decision nodes follow the interface:

```

1 class DNODE:
2     # children : dictionary[A] -> DNODE
3     def initialise(state ( $s_t$ ), depth ( $t$ ), context)
4     def select_action(context)
5     def backup(trial_return ( $R_t$ ), context)

```

And chance nodes:

```

1 class CNODE:
2     # children : dictionary[S] -> DNODE
3     def initialise(state ( $s_t$ ), action ( $a_t$ ), depth ( $t$ ), context)
4     def sample_outcome(context)
5     def backup(trial_return ( $R_t$ ), context)

```

The `run_trial` function can be written as:

```

1 def run_trial:
2     # root_node : DNODE
3     # mcts_mode : bool
4     t = 0
5     state = root_node.state
6     while (not selection_phase_ended(t, mcts_mode)):
7
8 def selection_phase_ended(t, mcts_mode):
9     if

```

urgh BRAIN POOP

TODO - copy the descriptions from DENTS, and adapt and add the psuedocode

## 2.3 Multi-Objective Reinforcement Learning

TODO: list

- MOMDP definition
- (Expected) utility
- Define an interface for pareto front and convex hull objects
- Define CHVI
- Should talk about multi-objective and/or contextual multi-armed bandits here?
- I'm planning on aligning this section with the recent MORL survey [1]
- Mention some deep MORL stuff, say that this work (given AlphaZero) is adjacent work

TODO: Follow CHMCTS and <https://arxiv.org/abs/2103.09568>

In this thesis we will follow a utility based approach to Multi-Objective Reinforcement learning similar to TODO: cite. For a full review of Multi-Objective Reinforcement Learning see TODO: cite. In this work we will specifically consider *linear utility* functions and the decision support scenario, which will be defined more precisely below.

This section defines the multi-objective counterparts to various definitions found in TODO: ref.

To specify problems with multiple objectives, the reward function of an MDP now outputs a vector of rewards, rather than a scalar reward.

**Definition 2.3.1.** A Multi-Objective Markov Decision Process (*MOMDP*) is a tuple  $\mathcal{M} = (\mathcal{S}, s_0, \mathcal{A}, \mathbf{R}, p, H)$ , where  $\mathcal{S}$  is a set of states,  $s_0 \in \mathcal{S}$  is an initial state,  $\mathcal{A}$  is a set of actions,  $\mathbf{V}(s, a)$  is a reward function  $\mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$ , where  $d$  is the dimension of the rewards and the MOMDP,  $p(s'|s, a)$  is a next state transition distribution  $\mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  and  $H \in \mathbb{N}$  is a finite-horizon time bound.

Now multi-objective trajectories are defined. Outside of this section, we may drop the prefix “multi-objective” where it should be clear from context, however we will continue to use bold typeface to denote any vector variables or functions.

**Definition 2.3.2.** *A multi-objective trajectory, is a sequence of state, action and rewards, that is induced by a policy  $\pi$  and MOMDP  $\mathcal{M}$  pair. Let the trials/-trajectory be  $\tau = (s_0, a_0, \mathbf{r}_0, s_1, a_1, \mathbf{r}_1, \dots, s_{H-1}, a_{H-1}, \mathbf{r}_{H-1}, s_H)$ , where  $a_t \sim \pi(\cdot|s_t)$ ,  $\mathbf{r}_t = \mathbf{R}(s_t, a_t)$  and  $s_{t+1} \sim \text{Succ}(s_t, a_t)$ . Notationally, we will write  $\tau \sim \pi$  to denote a sampled trial/trajectory with respect to a policy, where the MOMDP is implicit.*

*Sometimes it will be necessary to reason about trajectories with a horizon  $h < H$ , which will be denoted  $\tau_{:h} = (s_0, a_0, \mathbf{r}_0, s_1, a_1, \mathbf{r}_1, \dots, s_{h-1}, a_{h-1}, \mathbf{r}_{h-1}, s_h)$ .*

Similarly, multi-objective variants of the value of a policy needs to be defined:

**Definition 2.3.3.** *The value of a policy  $\pi$  from state  $s$  at time  $t$  is:*

$$\mathbf{V}^\pi(s; t) = \mathbb{E}_{\tau \sim \pi} \left[ \sum_{i=t}^{H-1} \mathbf{r}_i \middle| s_t = s \right]. \quad (2.31)$$

*The Q-value of a policy  $\pi$ , from state  $s$ , with action  $a$ , at time  $t$  is:*

$$\mathbf{Q}^\pi(s, a; t) = R(s, a) + \mathbb{E}_{s' \sim \text{Succ}(s, a)} [\mathbf{V}^\pi(s'; t + 1)]. \quad (2.32)$$

Now, in the corresponding single-objective section [TODO: ref](#) we have reached the point where we defined the optimal (Q-)value functions and the objective of single-objective reinforcement learning, where the maximum is taken over all possible policies. However, in a multi-objective setting there is no longer a *total ordering* over values, and so there maybe be multiple vectors that could be “optimal”. To resolve this issue, a *utility function* or *scalarisation function* is used to map multi-objective values to scalars.

**Definition 2.3.4.** *The ( $d$ -dimensional) Simplex consists of the set of  $d$ -dimensional vectors, whose entries are non-negative and sum to one. We denoted this as  $\Delta^d = \{\mathbf{w} \in \mathbb{R}^d | w_i \geq 0, \sum_i w_i = 1\}$ .*

*The elements of the  $d$ -dimensional Simplex we will call weight vector (or sometimes context in this thesis), as they will be used to specify preferences over the  $d$  dimensions of the reward function.*

**Definition 2.3.5.** A utility function (or scalarisation function)  $u : \mathbb{R}^d \times \Delta^d \rightarrow \mathbb{R}$  is used to map from a multi-objective value  $\mathbf{v} \in \mathbb{R}^d$  and a weighting over the objectives  $\mathbf{w} \in \Delta^d$  to a scalar value. That is, according to the utility function  $u(\cdot; \mathbf{w})$  the multi-objective value  $\mathbf{v}$  is mapped to the scalar value  $u(\mathbf{v}; \mathbf{w})$ .

TODO: comment here or in literature review about there being more types of scalarisation function that arent necessarily weighted by a weight, and ESR vs SER stuff

Of particular interest in this thesis is the *linear utility function* where the scalar value takes the form of a dot-product.

**Definition 2.3.6.** The linear utility function  $u_{\text{lin}}$  is the utility function defined by:

$$u_{\text{lin}}(\mathbf{v}; \mathbf{w}) = \mathbf{w}^\top \mathbf{v}. \quad (2.33)$$

Equiped with a scalarisation function and a weight vector any set of multi-objective values can be ordered. Letting  $\Pi$  be the set of all possible policies we can now define solution sets.

**Definition 2.3.7.** The undominated set of policies  $U(\Pi; u) \subseteq \Pi$ , with respect to a utility function  $u$ , is the set of policies for which there is a weight vector  $\mathbf{w} \in \Delta^d$  where the scalarised value is maximised:

$$U(\Pi; u) = \left\{ \pi \in \Pi \mid \exists \mathbf{w} \in \Delta^d. \forall \pi' \in \Pi : u(\mathbf{V}^\pi(s_0; 0); \mathbf{w}) \geq u(\mathbf{V}^{\pi'}(s_0; 0); \mathbf{w}) \right\}. \quad (2.34)$$

In particular, the convex hull of policies  $CH(\Pi)$  is the undominated set with respect to the linear utility function  $u_{\text{lin}}$ . That is  $CH(\Pi) = U(\Pi; u_{\text{lin}})$ .

TODO: write something here?

TODO: change d to D for the dimension, going to use d for a distance function

**Definition 2.3.8.** A set  $CS(\Pi; u) \subseteq \Pi$ , is a coverage set with respect to a utility function  $u$ , if for every weight vector  $\mathbf{w} \in \Delta^d$ , there is a policy  $\pi \in CS(\Pi; u)$  that

maximises the value of  $u(\cdot; \mathbf{w})$ . That is, for  $CS(\Pi; u)$  to be a coverage set, the following statement must be true:

$$\forall \mathbf{w} \in \Delta^d. \exists \pi \in CS(\Pi; u). \forall \pi' \in \Pi : u(\mathbf{V}^\pi(s_0; 0); \mathbf{w}) \geq u(\mathbf{V}^{\pi'}(s_0; 0); \mathbf{w}). \quad (2.35)$$

It can be shown that any  $CS(\Pi; u) \subseteq U(\Pi; u)$  *TODO: cite?*.

Again, in particular, any set  $CCS(\Pi)$  is a convex coverage set if it is a coverage set with respect to the linear utility function  $u_{\text{lin}}$ .

*TODO: some comments about undominated sets often being an infinite set of policies, but coverage sets often being finite and more feasible to compute*

*TODO: some comment about often computing the value set:*

**Definition 2.3.9.** The (multi-objective) value set with respect to a set of policies  $\Pi' \subseteq \Pi$  is defined as:

$$\mathcal{V}(\Pi') = \{\mathbf{V}^\pi(s_0; 0) | \pi \in \Pi'\}. \quad (2.36)$$

*TODO: Should acknowledge some things. Often we actually compute the value set of a convex coverage set. Often we compute a very specific convex coverage set, which is the extreme points of the convex hull. Also that the term convex hull is typically used to refer to any of the previous three sets (value set, convex coverage set, convex hull). And finally, say that often with methods that compute the value set can often use tagging to compute the policies after the fact, and cite some of the pomdp algorithms from LPK that actually explain the tagging*

Additionally, this thesis focusses on the decision support scenario as outlined in *TODO: cite*, with a linear utility function. In the decision support scenario the true weight vector is unknown, and so the objective is to compute a convex coverage set. When a convex coverage set is produced, it is then provided to a user that picks their most preferred policy or value from the coverage set. After this policy is selected, it can be used as a single solution to the problem that the user was trying to solve.

Moreover, in the case of MCTS algorithms, by having the user select a preferred policy, it implicitly forces the user to choose a preference over the objectives, as the

policy corresponds to a weight vector that it is optimal for. As MCTS algorithms are often used in an online fashion, where planning is interleaved with execution, this implicitly selected weight can be used for any online execution needed, effectively reducing the multi-objective problem into a single-objective problem.

TODO: make the above couple paragraphs not read like poo, maybe chatgpt it

### 2.3.1 Convex Hull Value Iteration

TODO: label for subsection]

TODO: write this

## 2.4 Multi-Objective Monte Carlo Tree Search

TODO: I think this whole section can just go in litrev

TODO: list

- Define the old methods (using the CH object methods, so clear that not doing direct arithmetic)
- Mention that old method could be written using the arithmetic of CHMCTS (but they don't)
- TODO: write about & make sure its implemented - its because just updating for 1 is more efficient in deterministic, and say that the additions can be implemented as updating for 1 value when deterministic
- Different flavours copy UCT action selection, but with different variants
- Link back to contributions and front load our results showing that all of the old methods don't explore correctly

TODO: There has been some prior work in multi-objective MCTS which we will outline here

TODO: Write out implementations of prior works using THTS

## 2.5 Sampling From Catagorical Distributions

TODO: list

- Talk about the alias method here
- Reference to chapter 4 section where talk about using this with THTS

TODO: clean this up generally, wrote it in a rush. Also trying not to use notation that I may want to use later. Would like

Much of the work in this thesis will involve sampling from catagorical distributions. Let  $f : \{1, \dots, m\} \rightarrow \mathbb{R}$  be the probability mass function of a catagorical distribution with  $m$  categories. Suppose that we want to sample  $i \sim f$ . A naive method to sample from  $f$  will take  $O(m)$  time, where a value is sampled from  $\text{Uniform}(0, 1)$  is often used as follows:

```

1 def sample_catagorical(f):
2     threshold ~ Uniform(0, 1)
3     i = 0
4     accumulated_mass = 0
5     while (accumulated_mass < threshold):
6         i += 1
7         accumulated_mass += f(i)
8     return i

```

However, the *Alias method* TODO: cite1, cite2 can instead be used, with  $O(m)$  preprocessing time to construct an *Alias table*, and can sample from  $f$  in  $O(1)$  time. In Figure TODO: ref we provide an example of an alias table. A value can be sampled using the alias table by sampling two random numbers, one from  $\text{Uniform}(\{1, \dots, m\})$  and one from  $\text{Uniform}(0, 1)$ . To sample from the alias table, one of the entries is sampled uniformly randomly using the sample from  $\text{Uniform}(\{1, \dots, m\})$ , each entry in the table contains three values, `threshold`, `cat_one` and `cat_two`, from which if we let  $a \sim \text{Uniform}(0, 1)$ , we would then return `cat_one` if  $a < \text{threshold}$  and `cat_two` otherwise. Psuedocode for this is as follows:

```

1 def sample_from_alias_table(alias_table):
2     index ~ Uniform({1, ..., m})

```

```
3 | cat_one, cat_two, threshold = alias_table[index]
4 | a ~ Uniform(0,1)
5 | if (a < threshold):
6 |     return cat_one
7 | return cat_two
```

In `TODO: cite` it is shown `TODO: check that it's proven` that an alias table can be constructed from an arbitrary probability mass function for a catagorical distribution, such that the probability of sampling any catagory from the alias table is identical to the probability of sampling it from the original probability mass function.

`TODO: verify that the example alias table maintains the correct masses for each category`

Following `TODO: cite`, we can construct an alias table as follows:

```
1 | def build_alias_table(f):
2 |     pass
```

`TODO: write the build alias table psuedocode (use wiki it was good)`



# 3

## Literature Review

### Contents

---

<b>3.1</b>	<b>Multi-Armed Bandits . . . . .</b>	<b>31</b>
<b>3.2</b>	<b>Reinforcement Learning . . . . .</b>	<b>32</b>
<b>3.3</b>	<b>Trial-Based Heuristic Tree Search and Monte-Carlo Tree Search . . . . .</b>	<b>32</b>
3.3.1	Trial Based Heuristic Tree Search . . . . .	32
3.3.2	Monte-Carlo Tree Search . . . . .	32
3.3.3	Maximum Entropy Tree Search . . . . .	32
<b>3.4</b>	<b>Multi-Objective Reinforcement Learning . . . . .</b>	<b>32</b>
<b>3.5</b>	<b>Multi-Objective Monte Carlo Tree Search . . . . .</b>	<b>33</b>
<b>3.6</b>	<b>Sampling Random Variables . . . . .</b>	<b>33</b>

---

TODO: currently this is a copy and paste of what I originally wrote for background chapter 2. Deleted parts which are irrelevant for litreview here (and vice versa for the background section).

TODO: I'm also going to use this as a space to paste papers I should write about as they come up while writing later chapters

### 3.1 Multi-Armed Bandits

TODO: Maybe dont need to cover this in litrev, but should talk about exploring bandits, UCT and contextual bandits either in background or in litrev

## 3.2 Reinforcement Learning

TODO: Intro should say that look at Sutton and Barto and something else for deep RL, for a more complete overview. Here we will just discuss papers that consider entropy in their work, as thats the most relevant part for this thesis.

TODO: list

- Talk about entropy and some of that work (probably a subsection)

## 3.3 Trial-Based Heuristic Tree Search and Monte-Carlo Tree Search

### 3.3.1 Trial Based Heuristic Tree Search

TODO: THTS paper

### 3.3.2 Monte-Carlo Tree Search

TODO: list

- Talk about the things that are ambiguous from literature (e.g. people will just say UCT, which originally presented doesn't run in `mcts_mode`, but often assumed it does)
- Should talk about multi-armed bandits here?

### 3.3.3 Maximum Entropy Tree Search

TODO: MENTS

## 3.4 Multi-Objective Reinforcement Learning

TODO: list

- Should talk about multi-objective and/or contextual multi-armed bandits here?
- Bunch of the work covered in recent MORL survey [1]

- Mention some deep MORL stuff, say that this work (given AlphaZero) is adjacent work

## 3.5 Multi-Objective Monte Carlo Tree Search

TODO: I think this whole section can just go in litrev

TODO: list

- Define the old methods (using the CH object methods, so clear that not doing direct arithmetic)
- Mention that old method could be written using the arithmetic of CHMCTS (but they don't)
- TODO: write about & make sure its implemented - its because just updating for 1 is more efficient in deterministic, and say that the additions can be implemented as updating for 1 value when deterministic
- Different flavours copy UCT action selection, but with different variants
- Link back to contributions and front load our results showing that all of the old methods don't explore correctly

## 3.6 Sampling Random Variables

TODO: Just citing the alias method papers?



# 4

## Monte Carlo Tree Search With Boltzmann Exploration

### Contents

---

4.1	Introduction . . . . .	35
4.2	Boltzmann Search . . . . .	36
4.3	Toy Environments . . . . .	36
4.4	Theoretical Results . . . . .	36
4.5	Empirical Results . . . . .	36
4.6	Full Results . . . . .	37

---

### 4.1 Introduction

TODO: list

- high level overview of DENTS work
- discuss how DENTS answers the research questions from introduction chapter
- state clearly that we're in single objective land here
- Comment about work exploring multi-armed bandits motivating this work

## 4.2 Boltzmann Search

TODO: list

- Recall MENTS
- Define BTS using THTS functions
- Define DENTS using THTS functions
- Discuss alias method variant (and complexity analysis) in a subsection?

## 4.3 Toy Environments

TODO: list

- Define D-chain stuff from the paper
- Define the D-chain with entropy trap
- Front load some results still

## 4.4 Theoretical Results

TODO: list

- add theoretical results

## 4.5 Empirical Results

TODO: list

- DChain
- GridWorlds
- Go

## 4.6 Full Results

TODO: there's a lot of figures for the D-chain environment, work out how to best fit them in? Or put them in this seperate section?





# 5

## Convex Hull Monte Carlo Tree Search

### Contents

---

5.1	Introduction . . . . .	39
5.2	Contextual Tree Search . . . . .	39
5.3	Contextual Zooming for Trees . . . . .	40
5.4	Convex Hull Monte Carlo Tree Search . . . . .	40
5.5	Results . . . . .	40

---

### 5.1 Introduction

TODO: list

- high level overview of CHMCTS work
- discuss how CHMCTS answers the research questions from introduction chapter
- moving into multi-objective land now
- Comment about CHVI and prior MOMCTS work motivating this

### 5.2 Contextual Tree Search

TODO: list

- Discuss need for context when doing multi-objective tree Search
  - Use an example env where left gives (1,0) and right gives (0,1), optimal policy picks just left or just right, but hypervolume based methods wont
  - Use previous work on these examples and show they dont do well bad
- Discuss how UCT = running a non-stationary UCB at each node, so given above discussion, there is work in contextual MAB
- Introduce contextual regret here

## 5.3 Contextual Zooming for Trees

TODO: list

- Give contextual zooming for trees algorithm
- Discussion on the contextual MAB to non-stationary contextual MAB stuff (CZT is to CZ what UCT is to UCB) (and what theory carry over)

## 5.4 Convex Hull Monte Carlo Tree Search

TODO: list

- Give convex hull monte carlo tree search
- Contextual zooming with the convex hull backups

## 5.5 Results

TODO: list

- Results from CHMCTS paper
- Get same plots from C++ code, but compare expected utility, rather than the confusing hypervolume ratio stuff

# 6

## Simplex Maps for Multi-Objective Monte Carlo Tree Search

### Contents

---

6.1	Introduction . . . . .	41
6.2	Simplex Maps . . . . .	42
6.3	Simplex Maps in Tree Search . . . . .	42
6.4	Theoretical Results . . . . .	42
6.5	Empirical Results . . . . .	42

---

### 6.1 Introduction

TODO: list

- high level overview of simplex maps work
- discuss how simplex maps answer the research questions from introduction chapter
- staying in multi-objective land now
- Motivated by CHMCTS being slow

## 6.2 Simplex Maps

TODO: list

- Define simplex map interface
- Give details on how to efficiently implement the interface with tree structures
- (Good diagram is everything here I think)

## 6.3 Simplex Maps in Tree Search

TODO: list

- Come up with better title for section
- Use simplex maps interface to create algorithms from the dents work
- Give a high level idea of what  $\delta$  parameter is (used in theory section)

## 6.4 Theoretical Results

TODO: list

- Convergence can build ontop of DENTS results
- Runtime bounds (better than  $O(2^D)$  which is what using convex hulls has)
- Simplex map has a diameter  $\delta$  (i.e. the furthest away a new context could be from a point in the map)
- Bounds can then come from that diameter (which is a parameter of the simplex map/algorithm) and DENTS results

## 6.5 Empirical Results

TODO: list

- Results from MO-Gymnasium
- Compare algorithms using expected utility

# 7

## Conclusion

### Contents

---

7.1	Summary of Contributions . . . . .	43
7.2	Future Work . . . . .	43

---

TODO: Something about we'll conclude by looking back at contributions and possible future work.

### 7.1 Summary of Contributions

TODO: go through each of the research questions and contributions, and write about how the work answers the research questions

### 7.2 Future Work

TODO: outline some avenues of potential future work



# Appendices







## List Of Appendices To Consider

- Multi Armed Bandits, maybe
- MMaybe from of the things in background are more appropriate as appendices?



# Bibliography

- [1] Conor F Hayes, Roxana Rădulescu, Eugenio Bargiacchi, Johan Källström, Matthew Macfarlane, Mathieu Reymond, Timothy Verstraeten, Luisa M Zintgraf, Richard Dazeley, Fredrik Heintz, et al. A practical guide to multi-objective reinforcement learning and planning. *Autonomous Agents and Multi-Agent Systems*, 36(1):26, 2022.
- [2] Thomas Keller and Malte Helmert. Trial-based heuristic tree search for finite horizon mdps. In *Proceedings of the International Conference on Automated Planning and Scheduling*, volume 23, pages 135–143, 2013.
- [3] Michael Painter. THTS++, <https://github.com/MWPainter/thts-plus-plus>.