

Possibility Theory for Reinforcement Learning

Tejas Gupta

Submitted as part of the honours requirements

Supervisor: Dr. Jeremie Houssineau

Division of Mathematical Sciences
School of Physical and Mathematical Sciences
Nanyang Technological University

April 2025

Abstract

A concise summary of why possibility theory is valuable in reinforcement learning, outlining the three major approaches and their key outcomes.

Acknowledgements

(Optional) Thank your supervisor, friends, colleagues.

Contents

Abstract	1
Acknowledgements	2
1 Introduction	5
1.1 Motivation and Context	5
1.2 Overview of Proposed Methods	5
1.3 Main Contributions	5
2 Background	6
2.1 Possibility Theory	6
2.1.1 Fuzzy Sets and Possibility Distributions	6
2.1.2 Additivity and Maxitivity	7
2.1.3 Normalization	7
2.1.4 Intersections and Unions: Conjunctions and Disjunctions	8
2.1.5 Fuzzy Measures and Integrals	8
2.2 Reinforcement Learning	8
2.2.1 Markov Decision Process	9
2.2.2 Deep-Q-Learning (DQN)	11
2.2.3 Actor-Critic Methods	12
2.2.4 Model-Based Reinforcement Learning	13
2.3 Possibility Theory and Reinforcement Learning	14
2.3.1 Distributional Reinforcement Learning	14
2.3.2 Possibilistic Q Learning	15
3 Possibilistic Q Values	16
3.1 Mean-Variance Networks	16
3.1.1 Possibilistic Bellman Equation	17
3.1.2 Loss Function	17
3.1.3 Action Selection Methods	18
3.2 Atomic Q Values	18
4 Proposed Approaches	20
4.1 Possibilistic Ensemble Q-Network	20
4.2 Model-Based MaxMax Possibility	20
5 Experimental Setup	21

5.1	Environments	21
5.2	Implementation Details	21
6	Results and Discussion	22
6.1	Performance Comparison	22
6.2	Insights	22
6.3	Limitations	22
7	Conclusion	23
7.1	Summary	23
7.2	Future Work	23
	References	23
A	Extra Details	25

Chapter 1

Introduction

1.1 Motivation and Context

1.2 Overview of Proposed Methods

1.3 Main Contributions

Chapter 2

Background

2.1 Possibility Theory

Possibility theory, introduced in Zadeh (1999), is a counterpart to probability theory that provides an alternative, flexible method of measuring and accounting for uncertainty. In this framework the uncertainty of an event is quantified by a possibility measure, which offers an alternative to model uncertainty due to incomplete knowledge. The possibility of an event can range from 0 to 1, where a value of 0 implies that the event is completely impossible and a value of 1 implies that the event is fully possible. In other words, possibility refers to the degree with which an event is possible given our current knowledge. This is in contrast to probability measures, where a probability of 1 implies that an event is statistically certain (or highly frequent), while a probability of 0.8 typically implies that the event happens with a frequency of 80%.

2.1.1 Fuzzy Sets and Possibility Distributions

Possibility theory was introduced as an extension to fuzzy sets in Zadeh (1999). A fuzzy set \tilde{A} is defined as a set of ordered pairs:

$$\tilde{A} = \{(x, \mu_{\tilde{A}}(x)) \mid x \in X\},$$

where $\mu_{\tilde{A}} : X \rightarrow [0, 1]$ is the membership function to the fuzzy set. The membership function over the set can also be understood as a *Possibility Distribution* $\hat{\pi}(x)$ over the set X . Analogous to probability theory, where the sum of the probabilities of all outcome states must be 1, a possibility distribution must ensure that at least one state is fully possible, i.e.,

$$\sup_{x \in X} \hat{\pi}(x) = 1.$$

The induced possibility measure for any subset $A \subseteq X$ is defined as the maximal value of $\hat{\pi}$ over the states:

$$\hat{\Pi}(A) = \sup\{\hat{\pi}(x) \mid x \in A\}.$$

This further implies that the union of two events is maxitive:

$$\hat{\Pi}(A \cup B) = \max\{\hat{\Pi}(A), \hat{\Pi}(B)\}.$$

Note that this holds even if A and B are not disjoint. In contrast to probability measures, where the probability of the union of disjoint events is the sum of the probabilities, possibility measures satisfy

$$\hat{\Pi}(\Omega) = 1, \quad \hat{\Pi}(\emptyset) = 1.$$

Dubois and Prade (2001) also introduced the notion of necessity, the dual of possibility, defined as

$$N(A) = \min\{1 - \hat{\pi}(x) \mid x \in A\} = 1 - \hat{\Pi}(\neg A).$$

Necessity quantifies the lack of plausibility of the complement of an event, so that possibility and necessity together can be interpreted as upper and lower probability bounds of imprecise probabilities (Dubois and Prade (1992)).

2.1.2 Additivity and Maxitivity

Probability measures are *additive*. For any two disjoint events A and B (i.e., $A \cap B = \emptyset$), the probability of their union is given by

$$P(A \cup B) = P(A) + P(B).$$

This additive property reflects the quantitative nature of probability, where the total weight is distributed among all outcomes.

In contrast, possibility measures are *maxitive* (or supremum-preserving) (Dubois and Prade (2007)). For any events A and B , the possibility measure of their union is given by

$$\hat{\Pi}(A \cup B) = \max\{\hat{\Pi}(A), \hat{\Pi}(B)\}.$$

This property implies that if at least one event is highly possible, then their union is considered highly possible. It allows possibility theory to express complete ignorance by simply assigning a possibility of 1 to all outcomes without forcing a partition of numerical weights.

2.1.3 Normalization

A probability distribution over an outcome space X requires that the probabilities of all states sum to 1:

$$\sum_{x \in X} P(x) = 1.$$

Even in situations of complete ignorance, a uniform distribution is imposed, which still assigns fractional probabilities to each outcome.

A possibility distribution, on the other hand, is normalized by requiring that at least one outcome has the maximal possibility:

$$\sup_{x \in X} \hat{\pi}(x) = 1.$$

This normalization permits complete ignorance to be represented trivially by assigning $\hat{\pi}(x) = 1$ for every x in X . Under such a distribution, each event has a necessity of 0, since

$$N(A) = 1 - \hat{\pi}(A^c) = 0,$$

when nothing is ruled out. This flexibility makes it easier to represent uncertainty qualitatively without imposing precise quantitative values.

2.1.4 Intersections and Unions: Conjunctions and Disjunctions

For independent events A and B , the probability of the joint event (the intersection) is typically given by the product:

$$P(A \cap B) = P(A) \cdot P(B).$$

Similarly, disjoint events have probabilities that add:

$$P(A \cup B) = P(A) + P(B).$$

In contrast, possibility theory uses triangular norms (t-norms) to model the logical AND (conjunction) of events (DUBOIS and (1982)). A common t-norm is the minimum operator, so that for events A and B with possibility distributions $\pi_A(x)$ and $\pi_B(x)$ respectively, the possibility distribution for the intersection is given by:

$$\pi_{A \cap B}(x) = \min\{\pi_A(x), \pi_B(x)\}.$$

This indicates that the possibility of a state satisfying both A and B is determined by the lesser possibility of the two. Dually, t-conorms (such as the maximum operator) are used for logical OR (disjunction):

$$\pi_{A \cup B}(x) = \max\{\pi_A(x), \pi_B(x)\}.$$

Thus, while probability theory uses multiplication (for independent events) and addition (for disjoint events), possibility theory replaces these operations with the minimum and maximum operators, respectively. This results in a very different arithmetic of uncertainty, which can simplify the handling of incomplete information.

2.1.5 Fuzzy Measures and Integrals

Possibility measures are a subset of fuzzy measures, which generalize classical measures by relaxing the requirement of additivity and requiring only monotonicity:

$$A \subseteq B \implies m(A) \leq m(B).$$

In possibility theory, rather than using the expected value computed via the Lebesgue integral, it is possible to aggregate outcomes using the Sugeno integral—a nonlinear operator based on the max and min operations. The Sugeno integral serves as an analogue to the Lebesgue integral and is particularly useful in qualitative decision-making scenarios where precise numeric integration is neither possible nor desired (Dubois and Prade (2015)).

2.2 Reinforcement Learning

Reinforcement Learning is a machine learning framework for an agent's sequential decision making in an environment. At each timestep, the agent observes the state in which it currently is, takes an action which moves it to another state, and collects a reward (the reward collected can be zero).

The notion of Actions, States, Rewards, and the associated stochastic transitions is formally known as the Markov Decision Process (MDP). Here we will discuss some core Reinforcement Learning concepts along with previous work employing possibility theory.

2.2.1 Markov Decision Process

A MDP is defined by the mathematical tuple (S, A, P, R, γ) where

- **State Space S :** refers to all possible states in an environment.
- **Action Space A_s :** refers to all possible actions available to the agent in the state s . In some formulations, the action space A might be the same across states.
- **Transition Probabilities $P(s' | s, a)$:** refers to the probability of transitioning to state s' by taking the action a in state s . These transitions can be either stochastic or deterministic.
- **Reward function $R(s, a, s')$:** is the immediate reward received by taking the action a in state s and transitioning to state s' . $R(s, a)$ refers to the expected reward received by taking action a in state s .
- **Discount Factor γ :** is the discounting factor of future rewards to determine the current value of the current state. A reward of 1 obtained after K steps is worth γ^K at the current step. Trivially, if γ is 1 then there is no discounting of future rewards.

As the name suggests, the Markov decision process also satisfies the Markov Property, i.e., the next state s' and the reward r only depend on the current state-action pair (s, a) ; all prior history is irrelevant.

The agent's behaviour in a state is characterised by its policy π , where $\pi(a | s)$ refers to the probability of the agent enacting a at state s . The goal of reinforcement learning is to find an optimal policy π^* that maximises cumulative rewards in an MDP.

R_t refers to the random variable denoting the reward the agent receives at timestep t . We can further define the cumulative rewards from the time step t as

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

Here, both the reward random variable and the cumulative reward random variable depend on the state at the current time t and the policy of the agent π . The expected cumulative reward under a given policy is represented by the state-value function $V^\pi(s)$.

$$V^\pi(s) = \mathbb{E}_\pi[G_t | S_t = s]$$

Similarly, an action value function (Q-value) $Q^\pi(s, a)$ can be defined as the expected cumulative return from state s if the agent takes action a .

$$Q^\pi(s) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a]$$

These expectations quantify how good an action or state is in terms of its expected cumulative rewards. Correspondingly, two policies can be compared on a given state by comparing the

value functions induced by that policy in that state. An optimal policy, hence, is the policy π^* that induces the optimal value function $V^*(s) = \max_{\pi} V^{\pi}(s)$ and $Q^*(s, a) = \max_{\pi} Q^{\pi}(s, a)$ for all s, a .

These expected values are also dependent on each other.

$$\begin{aligned} V^{\pi}(s) &= \mathbb{E}_{\pi}[Q^{\pi}(s, A) \mid S = s] \\ Q^{\pi}(s, a) &= \mathbb{E}^{\pi}[V^{\pi}(S') \mid S = s, A = a] \end{aligned}$$

By substituting the values further, one can construct a recursive relationship; this is also known as the Bellman Equation.

$$V^{\pi}(s) = \mathbb{E}_{\pi}[R_{t+1} + \gamma V^{\pi}(S_{t+1}) \mid S_t = s]$$

The state value of the current state is just the same as the transition reward and the discounted state value of the next state. A similar relationship exists for the action value function as follows

$$Q^{\pi}(s, a) = \mathbb{E}_{\pi}[R_{t+1} + \gamma Q^{\pi}(S'_{t+1}, A_{t+1}) \mid S_t = s]$$

The definition of the recursive expectations can be fully expanded as follows:

$$\begin{aligned} V^{\pi}(s) &= \sum_{a \in A} \pi(a \mid s) \sum_{s' \in S} P(s' \mid s, a) [R(s, a, s') + \gamma V^{\pi}(s')] \\ Q^{\pi}(s, a) &= \sum_{s'} P(s' \mid s, a) [R(s, a, s') + \gamma \sum_{a'} \pi(a' \mid s'), Q^{\pi}(s', a')] \end{aligned}$$

For a given policy, the Bellman Equations are linear. However, for an optimal policy, we have nonlinear maximisation operations as follows:

$$V^*(s) = \max_{a \in A} \mathbb{E}[R_{t+1} + \gamma V^*(S_{t+1}) \mid S_t = s]$$

This gives the intuitive result that the optimal value of a state is the same as the expected value of taking the best action from the state. The same result also applies to the Q-function:

$$Q^*(s, a) = \mathbb{E}[R_{t+1} + \gamma \max_{a'} Q^*(S_{t+1}, a') \mid S_t = s, A_t = a]$$

In a finite state and action space, it is possible to solve the Bellman Optimality Equations to get the optimal values using value iteration or policy iteration, making it possible to calculate V^* from which it is trivial to deduce an optimal policy. However, in larger and continuous environments, this is no longer feasible. Thus, reinforcement learning algorithms attempt to either learn the value functions directly or learn a maximising policy by experiencing the MDP.

2.2.2 Deep-Q-Learning (DQN)

Q-Learning is a category of reinforcement learning algorithms that focus on learning the optimal Q^* value for each state action pair by iterative updates. In its simple tabular form, the iteration happens as follows:

$$Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$

Here α is the learning rate and (s, a, r, s') is one of the experienced transitions. This update is based on the Bellman Update discussed before. Watkins and Dayan (1992) showed that this tabular Q-Learning converges to the optimal value if all states are visited infinitely during the learning process and the learning rates satisfy:

$$\sum_{t=0}^{\infty} \alpha_t(s, a) = \infty$$

$$\sum_{t=0}^{\infty} \alpha_t^2(s, a) < \infty$$

As mentioned before, the tabular approach to Q-Learning is not feasible for continuous or large environments. Deep Q Learning provides an alternative method to functionally approximate the Q-Values using deep neural networks. The function $Q(s, a | \theta)$ is parametrised by weights θ . In general, it is possible to approximate the Q-Values using other methods (such as Linear Functions). The method for doing Deep Q Learning was first introduced Mnih et al. (2015) where the authors achieved human level performance in various ATARI games. In their paper, the network takes the state of the game (in the form of a raw image) and outputs the approximate Q-Values for each of the set of discrete actions. Common to other Deep Learning Method, Stochastic Gradient Descent is used to update the Q networks, with the loss function derived from the temporal difference (TD) error. Particularly, the Bellman Backup is used as the target for training, i.e the update target y for the transition (s, a, r, s') is:

$$y = r + \gamma \max_{a'} Q(s', a'; \theta^-)$$

where θ^- is the parameters of the target network. A target network is typically a lagged copy of the main Q-Network and helps stabilise the learning process. The target network can either be continuously updated using Polyak Averaging $\theta^- \leftarrow \theta^- + \alpha(\theta - \theta^-)$ or it can be copied from the main network periodically during the learning process. The loss function can be constructed to minimise the mean squared error between y and its own outputs:

$$L(\theta) = \mathbb{E}[(y - Q(s, a; \theta))^2]$$

The loss is calculated over a batch D which is a set of tuples (s, a, r, s') experienced by the agent. Taking the gradient of the loss:

$$\nabla_{\theta} L(\theta) = \mathbb{E}_{(s, a, r, s') \sim D} [2 \cdot (Q(s, a; \theta) - y) \cdot \nabla_{\theta} Q(s, a; \theta)]$$

The network is updated by moving the parameters θ to minimise the loss:

$$\theta \leftarrow \theta - \eta \cdot \nabla_{\theta} L(\theta)$$

In practice, the networks are not trained using sequential experiences as this can lead to divergence and instability in the learning process. Instead, Mnih et al. (2015) introduced experiential learning. The transitions (s, a, r, s') experienced during training are stored in a replay buffer. During the training step, a mini-batch is randomly sampled from the replay buffer. The random sampling breaks temporal correlations in the data and generally results in smoother learning. Reusing past transitions also improves data efficiency as each transition can be used multiple times to improve the learning process. Q-Learning is a type of off-policy learning, as the method does not directly learn the policy, rather it only learns an estimate of optimal Q-Value Q^* ; this means that the Q-Values can be trained from any observed transitions in a given environment.

2.2.3 Actor-Critic Methods

Actor Critic Methods are a type of On-Policy reinforcement learning algorithm, in that they explicitly learn the policy of the actor along with a critic (usually the value functions). In general, the actor decides what actions to take based on the state, while the critic estimates the advantage of the actions. The actor is updated to maximise the advantage estimate of the critic. Actor-Critic methods can handle continuous and larger action spaces better than Q-Learning, while also learning the value estimates to perform stable policy gradient updates.

In policy gradient methods, the policy $\pi_\theta(s)$ is parametrised by θ , this is often implemented using neural networks. An objective function $J(\theta)$ is defined as the expected return from the environment following policy parametrised by θ from some starting distribution of states. The policy gradient theorem provides a way to express the graient of the objective function $J(\theta)$ as

$$\nabla_\theta J(\theta) = \mathbb{E}_{s \sim d^\pi, a \sim \pi_\theta} [Q^\pi(s, a) \nabla_\theta \log \pi_\theta(a | s)]$$

where s is distributed by d^π (the visitation frequency of the states under the policy π_θ) and a is distributed by the policy π_θ . This means that the gradient of expected reward can be calculated as the sum of action values weighted by

$$\frac{\nabla_\theta \pi_\theta(a | s)}{\pi_\theta(a | s)}$$

The parameter can then updated using

$$\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$$

In practice, one uses the advantage function $A(s, a) = Q(s, a) - V(s)$ instead of the $Q(s, a)$ value directly as that increases stability in learning. It is shown that replacing the Q-Value with the advantage function does not introduce a bias.

The actor network is then updated with:

$$\Delta \theta_{\text{actor}} \propto \nabla_\theta \log \pi_\theta(a_t | s_t) \hat{A}_t$$

where \hat{a} is the empirical advantage.

The critic in actor critic environments is commonly the optimal state-value function approximate $V_w(s)$ or the optimal action-value approximation $Q_w(s, a)$. The critic can be trained

on similar to Q-Learning using temporal difference (TD) learning. The state-value critic, for example, can be updated by minimising the error δ_t :

$$\delta_t = r_{t+1} + \gamma V_{w'}(s_{t+1}) - V_w(s_t)$$

Similar to the DQN learning discussed before, there is usually a target critic and online critic to stabilise learning. Importantly, $r_{t+1} + \gamma V_{w'}(s_{t+1})$ is an estimator of the Q Value $Q(s, a)$ where the actor a is determined by the current policy π_θ ; therefore δ_t itself is an estimator for the advantage function. Hence, the actor's gradient update can be implemented as

$$\Delta\theta \propto \Delta_\theta \log \pi_\theta(a_t | s_t)$$

which nudges the policy to increase the probability of action a_t if the observed reward is greater than the expected value estimated by $V_w(s)$. In other words, the critic criticises the actions and the actor uses this to improve its policy.

It is also possible to have deterministic policies instead of stochastic policy described above. In the deterministic policy gradient approach (DPG) $a = \mu_\theta$ represents the deterministic policy produced by the actor. The update to the actors policy is implemented as:

$$\nabla_\theta J \approx \mathbb{E}_{s \sim D} [\nabla_a Q_w(s, a) | a = \mu_\theta(s), \nabla_\theta \mu_\theta(s)]$$

Here, the action-value function itself is the critic, and μ_θ is updated to maximise the Q-Value. The Q value is trained using the mini-batches, similar to in DQN.

2.2.4 Model-Based Reinforcement Learning

The algorithms discussed earlier are all model-free, in that, none of the methods maintain an internal "model" of the environment. In contrast, there exists a category of algorithms that use and/or learn a model of the environment. This can include learning the state transitions $\hat{P}(s' | s, a)$ and the reward function $\hat{R}(s, a)$. The agent is able to plan with the environment model to evaluate its actions without enacting them in the real environment.

In MDP, if the true transition and reward functions are known, it is possible to infer the optimal policies by dynamic programming, this is also referred to as planning. However, in most applications, the exact P, R are unknown and the various Model Based approaches attempt to learn a model of environment through repeated interaction. The learned P, R can then be used for planning in the environment. A major benefit of model based approaches is the improved sample efficiency during training as the agent is able to simulate experience using the environment model instead of taking real actions. Model Based learning can be particularly beneficial if taking samples in the real environment is costly or limited.

In practice, learning the model can be treated as a supervised learning task as it is possible to treat each recorded transition (s, a, r, s') as a training sample for the models \hat{P} and \hat{R} . The model can help both in action selection and with value iteration. Dyna-Q Learning is common approach for value iteration, where the Q values are updated by simulated planning steps that are intermixed with real experiences. The model can also help in action selection, for example, it is possible to do Monte Carlo Tree Search to simulate action sequences and select the best performing action.

One of the essential challenges of model based learning is learning an accurate model. A biased or imperfect model can create result in non-optimal policies. There are many different approaches to account for this including: shorter planning horizons, uncertainty bounds on the model and ensemble models.

2.3 Possibility Theory and Reinforcement Learning

Reinforcement Learning algorithm primarily deal with uncertainty in two different forms: aleatoric uncertainty and epistemic uncertainty. Aleatoric uncertainty refers to the uncertainty in the environment because of the inherent randomness; this could include randomness in state transitions and stochastic rewards. Epistemic uncertainty, on the other hand, refers to the uncertainty in the environment because of our lack of knowledge/ information. Possibility Theory can be helpful to evaluate and make use of the latter.

2.3.1 Distributional Reinforcement Learning

One possible method of handling uncertainty in the environment is by maintaining distributions over $Q(s, a)$ or $V(s, a)$ instead of a single value. Along the same line, distributional Reinforcement Learning, introduced in Bellemare et al. (2017), provides a novel method to handle uncertainty in the reward distribution by modelling a distribution Z^π over Q Values. In particular, instead of learning only

$$Q^\pi(s, a) = \mathbb{E}[Z^\pi(s, a)]$$

where Z^π is the random return, satisfying the above expectation and the recursive relationship

$$Z^\pi(s, a) \stackrel{D}{=} R(s, a) + \gamma Z^\pi(s', a')$$

.

The above return also satisfies the Bellman Equation:

$$Z^*(s, a) \stackrel{D}{=} R_{t+1} + \gamma Z^*\left(S_{t+1}, \arg \max_{a'} \mathbb{E}[Z^*(S_{t+1}, a')]\right)$$

The algorithm C51 in Bellemare et al. (2017) works by maintaining a probability distribution over 51 possible values $\{z_1, z_2, \dots, z_{51}\}$, which are placed uniformly over the range of returns. In particular, instead of the Q-Network outputting a single float value of state action pair, the network returns a probability tuple $(P(Z(s, a) = z_1), P(Z(s, a) = z_2), \dots, P(Z(s, a) = z_{51}))$. For training, a target distribution is calculated as:

$$T(s, a) \stackrel{D}{=} r + \gamma Z(s', a^*)$$

where a^* is the greedy action that maximises $\mathbb{E}[Z(s', a)]$.

One drawback of the approach is that the two different uncertainties cannot be easily decoupled. If information about the epistemic uncertainty of the distribution was available, it would be possible to incentive the agent to explore states with less certainty.

2.3.2 Possibilistic Q Learning

Possibilistic Q Learning, as introduced by Thomas and Houssineau (2025), extends traditional Q-Learning by explicitly accounting for both aleatoric and epistemic uncertainty through the use of possibility theory. In this framework, epistemic uncertainty is modeled with possibility functions, allowing for a clear separation between uncertainty due to inherent randomness (aleatoric) and uncertainty due to lack of knowledge (epistemic).

A key concept in this approach is the *maximum expected value*, defined as

$$\bar{\mathbb{E}}(\Phi(\psi)) = \sup_{\psi \in \Psi} \{ \Phi(\psi) f_{\psi}(\psi) \},$$

where f_{ψ} denotes the possibility function over the parameter ψ . This operator selects the highest weighted value of $\Phi(\psi)$, thereby capturing an optimistic estimate based on the current knowledge. Correspondingly, the most credible value of ψ is given by

$$\mathbb{E}^*(\psi) = \arg \max_{\psi \in \Psi} f_{\psi}(\psi).$$

Using these definitions, one can derive a recursive formulation for the maximum expected Q-value:

$$\bar{Q}(s, a) = \bar{\mathbb{E}}_{S'} \left[\bar{r}(s, a, S') + \max_{a'} \bar{Q}(s', a' | S') \right].$$

Here, $\bar{r}(s, a, S')$ represents the expected reward computed with the possibility-based model, and $\bar{Q}(s', a' | S')$ denotes the optimistic Q-value at the next state s' .

In this model, a possibility function $P(\cdot | s, a)$ is maintained over the state transitions. As training progresses, this possibility function is updated along with the maximum expected value, adapting the model's uncertainty estimates based on new observations. Since the initial Q-values in the tabular setting are set to an upper bound, the algorithm naturally favors exploration in states with high uncertainty. Over time, as more data is gathered, the optimistic estimates are refined toward the true Q-values.

Chapter 3

Possibilistic Q Values

Typically in Q-Learning based algorithm, the Q-Network outputs a tuple of scalars representing the networks estimates of expected cumulative reward of the particular action a in state s . Exploration, then, driven by external stochasticity (e.g., with ϵ -greedy strategies). This strategy is unable to account for the aleatoric uncertainties in the environment and the models own epistemic uncertainty in the estimated Q values. By representing uncertainty, it can be possible to guide the agent to explore actions with higher uncertainty to gather more experience (in the form of transitions (s, a, r, s')) and increase certainty in its Q values. Here we will present two different algorithms using possibility distributions over $Q(s, a)$.

3.1 Mean-Variance Networks

In this proposed possibilistic approach, instead of focusing on a single scalar for $Q(s, a)$, we evaluate a possibility distribution over the Q-values, denoted by $f(q \mid s, a)$. The distribution is parameterized by a mean and a variance; that is, the Q-network outputs two parameters for each state-action pair:

- **Mean** $\mu(s, a)$ — the expected value of the return (equivalent to $Q(s, a)$).
- **Variance** $\sigma^2(s, a)$ — which quantifies the uncertainty in the value of $\mu(s, a)$.

Together, these parameters define a Gaussian-like membership function for the state-action pair:

$$f(q \mid s, a) = \exp\left(-\frac{(q - \mu(s, a))^2}{2\sigma^2(s, a)}\right).$$

By definition, this membership function satisfies the normalization condition

$$\sup_q f(q \mid s, a) = 1,$$

which is achieved at $q = \mu(s, a)$. This aligns with the intuitive belief that the most credible value of the return is the mean.

3.1.1 Possibilistic Bellman Equation

A possibilistic Bellman equation can be formulated for the above. Given an observed transition (s, a, r, s') , we select the action

$$a' = \arg \max_{a'} \mu(s', a')$$

and define the target parameters as:

$$\begin{aligned} \mu_{\text{target}} &= r + \gamma \mu(s', a'), \\ \sigma_{\text{target}}^2 &= \gamma^2 \sigma^2(s', a'). \end{aligned}$$

This follows directly from the conventional Bellman equation for Q-learning,

$$Q(s, a) = r + \gamma \max_{a'} Q(s', a'),$$

and, in our case, we propagate both the mean and the uncertainty. Thus, the target distribution can be modeled as:

$$\mathcal{T}Q(s, a) \sim \mathcal{N}\left(r + \gamma \mu(s', a'), \gamma^2 \sigma^2(s', a')\right).$$

Note that the γ^2 factor in the variance indicates that, if the mean estimate is accurate, then over time the variance should decrease, reflecting an increase in certainty regarding the μ estimates.

3.1.2 Loss Function

To update the network, the discrepancy between the target distribution and the current estimated distribution must be minimized. Divergence metrics common in probability theory serve as proxies for a distance measure between our possibilistic distributions. In particular, we consider:

Kullback–Leibler Divergence

For two Gaussian distributions, the KL divergence is given by

$$D_{\text{KL}}\left(\mathcal{N}(\mu_1, \sigma_1^2) \parallel \mathcal{N}(\mu_2, \sigma_2^2)\right) = \frac{1}{2} \left[\log \left(\frac{\sigma_2^2}{\sigma_1^2} \right) + \frac{\sigma_1^2}{\sigma_2^2} + \frac{(\mu_1 - \mu_2)^2}{\sigma_2^2} - 1 \right],$$

which measures how "surprised" the current estimate $\mathcal{N}(\mu_1, \sigma_1^2)$ would be if the target were $\mathcal{N}(\mu_2, \sigma_2^2)$. A note about D_{KL} is in the Appendix A.

Wasserstein-2 Metric

The Wasserstein-2 metric (also known as the 2nd-order Wasserstein distance) between two Gaussian distributions is

$$W_2^2\left(\mathcal{N}(\mu_1, \sigma_1^2), \mathcal{N}(\mu_2, \sigma_2^2)\right) = (\mu_1 - \mu_2)^2 + (\sigma_1 - \sigma_2)^2.$$

This metric provides a geometric measure of the distance between the two distributions, incorporating both the differences in means and differences in spread, and in our case is equivalent to a mean-squared error between the network parameters.

3.1.3 Action Selection Methods

Information about uncertainty can be utilized to incentivize exploration. We consider two different methods of action selection.

Log-Variance-Weighted Selection

In this method, the action is chosen by combining the mean and the log-variance of the Q-value:

$$a^* = \arg \max_a \{ \mu(s, a) + \beta \cdot \log \sigma^2(s, a) \}.$$

Here, β is a hyperparameter that adjusts the influence of the uncertainty (measured by $\log \sigma^2(s, a)$) on the action selection. **Why is this a good way to choice, why not without the log**

Maximum Expected Value

Alternatively, we can utilize the notion of maximum expected value introduced in Thomas and Houssineau (2025). The optimistic estimate of the Q-value is defined as

$$\bar{Q}(s, a) = \sup_{q \in \mathbb{Q}} \{ q f(q \mid s, a) \}.$$

Under the assumption that $f(q \mid s, a)$ is a Gaussian possibility function, this maximum has a closed-form solution (see Lemma 1 in Appendix):

$$\bar{Q}(s, a) = \frac{\mu(s, a) + \sqrt{\mu(s, a)^2 + 4\sigma^2(s, a)}}{2}.$$

3.2 Atomic Q Values

As discussed before, Bellemare et al. (2017) introduced distributional RL, which parameterizes a distribution over Q-values using fixed atomic values and assigns probabilities to these fixed atoms. This approach can potentially be enhanced using possibility theory—that is, by generating possibility values for the canonical atomic returns.

Let $f(q \mid s, a)$ be the mapping that assigns to each $q \in \mathbb{Q}$ (where \mathbb{Q} is a finite ordered tuple of possible q values, e.g., $\{q_1, q_2, \dots, q_N\}$) the possibility that q is the true Q-value for the state–action pair (s, a) . An informed prior over the Q-values can be formed by setting

$$f(q \mid s, a) = 1 \quad \forall q \in \mathbb{Q}, s, a,$$

i.e., initially all values of q are considered fully possible.

For a transition (s, a, r, s') , we can form a likelihood function by using a Gaussian kernel. Specifically, the likelihood of observing a reward r given that the true return is q_j is defined as:

$$p_R(r \mid q_j) = \exp\left(-\frac{(q_j - r)^2}{2\sigma^2}\right).$$

This likelihood is used to update the possibility function $f(q \mid s, a)$ via a possibilistic version of Bayes' rule. The updated possibility of the atom q_j is given by:

$$f(q_j \mid s, a) \leftarrow \frac{p_R(r \mid q_j) f(q_j \mid s, a)}{\sup_{k \in \{1, 2, \dots, N\}} \{p_R(r \mid q_k) f(q_k \mid s, a)\}}.$$

Since $f(q_j \mid s, a)$ is initialized to 1 for all q_j , the first update simplifies to the ratio of likelihoods.

Note that by construction, for every s, a there exists at least one j such that

$$f(q_j \mid s, a) = 1,$$

satisfying the normalization condition in possibility theory.

In the discrete setting, the function $f(q_j \mid s, a)$ can be parameterized with a table. The optimistic Q-value can then be computed using the maximum expected value method:

$$\bar{Q}(s, a) = \sup_{j \in \{1, \dots, N\}} \{q_j f(q_j \mid s, a)\}.$$

This means that initially, when the possibility scores are all 1, the maximum Q-value is chosen (i.e., all actions appear highly promising). As learning proceeds, the credibility of overly optimistic q_j values will decrease and $\bar{Q}(s, a)$ will favor the Q-values with high credibility.

This framework can be generalized to the deep learning case by using a deep neural network as a function approximator for $f(q_j \mid s, a)$.

A Bellman target distribution can be formulated as

$$\mathcal{T}Q = r + \gamma f(s', a'),$$

where $a' = \arg \max_a \bar{Q}(s', a)$. **How to account for the shift here?** For the deep learning update, the possibility distribution must be renormalized as follows:

$$f'(q_j \mid s, a) = \frac{f(q_j \mid s, a)}{\sup_{q \in \mathbb{Q}} f(q \mid s, a)}.$$

A loss function such as the L_2 divergence loss can be used for gradient descent. For example, if $\hat{f}(q_j \mid s, a)$ denotes the target possibility distribution obtained via the Bellman backup, we can define the loss as:

$$L(s, a) = \sum_{j=1}^N \left(f(q_j \mid s, a) - \hat{f}(q_j \mid s, a) \right)^2.$$

Minimizing this loss updates the network parameters so that the predicted possibility distribution aligns with the target distribution.

Chapter 4

Proposed Approaches

4.1 Possibilistic Ensemble Q-Network

4.2 Model-Based MaxMax Possibility

Chapter 5

Experimental Setup

5.1 Environments

5.2 Implementation Details

Chapter 6

Results and Discussion

6.1 Performance Comparison

6.2 Insights

6.3 Limitations

Chapter 7

Conclusion

7.1 Summary

7.2 Future Work

Bibliography

- Bellemare, M. G., Dabney, W., and Munos, R. (2017). A distributional perspective on reinforcement learning.
- Cover, T. M. and Thomas, J. A. (1991). *Elements of Information Theory*. John Wiley & Sons.
- DUBOIS, D. and and, H. P. (1982). A class of fuzzy measures based on triangular norms a general framework for the combination of uncertain information. *International Journal of General Systems*, 8(1):43–61.
- Dubois, D. and Prade, H. (1992). When upper probabilities are possibility measures. *Fuzzy Sets and Systems*, 49(1):65–74.
- Dubois, D. and Prade, H. (2001). Possibility theory, probability theory and multiple-valued logics: A clarification. *Ann. Math. Artif. Intell.*, 32:35–66.
- Dubois, D. and Prade, H. (2007). Possibility theory. *Scholarpedia*, 2(10):2074. revision #137677.
- Dubois, D. and Prade, H. (2015). Possibility theory and its applications: Where do we stand? *Mathware and Soft Computing Magazine*, 18.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., and Hassabis, D. (2015). Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533.
- Thomas, J. and Houssineau, J. (2025). Possibilistic q-learning: Uncertainty modelling for tuning-free optimism. Unpublished manuscript.
- Watkins, C. J. C. H. and Dayan, P. (1992). Q-learning. *Machine Learning*, 8(3):279–292.
- Zadeh, L. (1999). Fuzzy sets as a basis for a theory of possibility. *Fuzzy Sets and Systems*, 100:9–34.

Appendix A

Extra Details

Lemma 1 (Closed Form Maximum of $q \cdot f(q)$ for Gaussian Possibility). *Let $f(q)$ be a Gaussian-shaped possibility function:*

$$f(q) = \exp\left(-\frac{(q - \mu)^2}{2\sigma^2}\right)$$

Define

$$g(q) = q \cdot f(q) = q \cdot \exp\left(-\frac{(q - \mu)^2}{2\sigma^2}\right).$$

Then the supremum of $g(q)$ over $q \in \mathbb{R}$ is achieved at

$$q^* = \frac{\mu + \sqrt{\mu^2 + 4\sigma^2}}{2}.$$

Proof. We differentiate $g(q)$ with respect to q :

$$g(q) = q \cdot \exp\left(-\frac{(q - \mu)^2}{2\sigma^2}\right),$$

$$g'(q) = \exp\left(-\frac{(q - \mu)^2}{2\sigma^2}\right) \left(1 - \frac{q(q - \mu)}{\sigma^2}\right).$$

Setting $g'(q) = 0$, we solve

$$1 - \frac{q(q - \mu)}{\sigma^2} = 0 \quad \Rightarrow \quad q^2 - \mu q - \sigma^2 = 0.$$

Solving this quadratic equation yields

$$q^* = \frac{\mu + \sqrt{\mu^2 + 4\sigma^2}}{2}.$$

This is the unique maximizer of $g(q)$ over \mathbb{R} , as $g''(q) < 0$ at this point. □

Kullback-Leibler Divergence

In Cover and Thomas (1991), the relative entropy (or Kullback-Leibler divergence) is introduced as the expected logarithm of the likelihood ratio. It measures the inefficiency of assuming that the distribution is q when the true distribution is p .

Discrete Case

For discrete probability mass functions $p(x)$ and $q(x)$, the KL divergence is defined as

$$D(p\|q) = \sum_x p(x) \log \frac{p(x)}{q(x)}.$$

Continuous Case

For continuous probability density functions $p(x)$ and $q(x)$, the KL divergence is defined as

$$D(p\|q) = \int p(x) \log \frac{p(x)}{q(x)} dx.$$

KL Divergence Between Two Gaussian Distributions

Consider two Gaussian distributions:

$$p(x) = \mathcal{N}(\mu_1, \sigma_1^2), \quad q(x) = \mathcal{N}(\mu_2, \sigma_2^2).$$

The KL divergence is given by

$$\begin{aligned} KL(p\|q) &= - \int p(x) \log q(x) dx + \int p(x) \log p(x) dx \\ &= \frac{1}{2} \log(2\pi\sigma_2^2) + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2} \left(1 + \log(2\pi\sigma_1^2) \right) \\ &= \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}. \end{aligned}$$