# 1 Introduction

In this document, we will explore a key equation in Bayesian decision theory and reinforcement learning:

$$P(a_t|b_{t-1}, g) = \frac{\exp(-\beta \hat{Q}_g(b_{t-1}, a))}{\sum_{a'} \exp(-\beta \hat{Q}_g(b_{t-1}, a'))}$$

To understand this equation and its implications, we'll delve into the basics of Bayesian analysis, its application to decision-making, and how it relates to reinforcement learning.

# 2 Basics of Bayesian Analysis

## 2.1 Bayes' Theorem

The foundation of Bayesian analysis is Bayes' theorem:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Where:

- $P(A|B)$ is the posterior probability of A given B

- $P(B|A)$ is the likelihood of B given A

- $P(A)$ is the prior probability of A

- $P(B)$ is the marginal likelihood of B

## 2.2 Prior and Posterior Distributions

In Bayesian analysis:

- The prior distribution represents our initial belief about a parameter before observing data.

- The posterior distribution updates our belief after observing data.

## 2.3 Likelihood

The likelihood function $P(B|A)$ represents the probability of observing the data given a particular parameter value.

# 3 Bayesian Decision Theory

Bayesian decision theory applies Bayesian inference to decision-making problems.

## 3.1 Expected Utility

In decision theory, we often work with expected utility:

$$EU(a) = \sum_s P(s|a)U(s)$$

Where:

- $EU(a)$ is the expected utility of action $a$
- $P(s|a)$ is the probability of state $s$ given action $a$
- $U(s)$ is the utility of state $s$

# 4 Analysis of the Key Equation

Now, let's analyze our key equation in detail:

$$P(a_t|b_{t-1}, g) = \frac{\exp(-\beta \hat{Q}_g(b_{t-1}, a))}{\sum_{a'} \exp(-\beta \hat{Q}_g(b_{t-1}, a'))}$$

## 4.1 Components of the Equation

1. $P(a_t|b_{t-1}, g)$:

   - This is the probability of taking action $a_t$ at time $t$.
   - It's conditional on $b_{t-1}$ (the belief state at time $t-1$) and $g$ (the goal or task).

2. $\exp(-\beta \hat{Q}_g(b_{t-1}, a))$:

   - $\hat{Q}_g(b_{t-1}, a)$ is an estimated Q-value function.
   - It represents the expected future reward of taking action $a$ in belief state $b_{t-1}$ for goal $g$.
   - $\beta$ is a temperature parameter that controls the randomness of the policy.
   - The negative sign inverts the Q-value, so lower Q-values result in higher probabilities.
   - exp() is the exponential function, which ensures all values are positive.

3. $\sum_{a'} \exp(-\beta \hat{Q}_g(b_{t-1}, a'))$:

   - This is the sum over all possible actions $a'$.
   - It serves as a normalization factor to ensure the probabilities sum to 1.

## 4.2 Interpretation

- Actions with higher Q-values (lower $-\beta\hat{Q}_g$) are assigned higher probabilities.

- The $\beta$ parameter controls the "sharpness" of the distribution:
  - As $\beta \to \infty$, it approaches a deterministic policy (always choose the best action).
  - As $\beta \to 0$, it approaches a uniform distribution (all actions equally likely).

- $P(a_t|b_{t-1}, g)$ gives us a posterior distribution over actions, updated based on our current belief and goal.

## 4.3 Bayesian Perspective

From a Bayesian viewpoint:

- $b_{t-1}$ represents our prior belief about the state of the world.

- $\hat{Q}_g$ can be seen as incorporating our likelihood model of how actions affect the world and our utility function for different outcomes.

- The equation gives us a posterior distribution over actions, updated based on our current belief and goal.

## 4.4 Relation to Softmax Function

The equation is a softmax function, which is related to the Boltzmann distribution:

$$P(x) = \frac{\exp(-\beta E(x))}{\sum_{x'} \exp(-\beta E(x'))}$$

In our case, $E(x)$ is replaced by $\hat{Q}_g(b_{t-1}, a)$.

# 5 Applications in Reinforcement Learning

In reinforcement learning:

- $\hat{Q}_g(b_{t-1}, a)$ is often learned through experience.

- This softmax policy allows for exploration (trying suboptimal actions) while still favoring actions with higher expected rewards.

- $\beta$ controls the exploration-exploitation trade-off.

# 6    Conclusion

The equation we've analyzed represents a probabilistic policy for action selection, grounded in Bayesian decision theory and commonly used in reinforcement learning. It balances the exploitation of known good actions with exploration of potentially better alternatives, all within a Bayesian framework of updating beliefs based on evidence and prior knowledge.

This form of the Boltzmann distribution, also known as the softmax function in machine learning, is a powerful tool for converting value estimates into action probabilities. It provides a flexible framework for decision-making under uncertainty, allowing for adaptive behavior in complex, dynamic environments.