

# Study Material - Youtube

## Document Information

- **Generated:** 2025-08-26 08:12:41
- **Source:** <https://www.youtube.com/watch?v=kCwdMhKO0x8>
- **Platform:** Youtube
- **Word Count:** 2,404 words
- **Estimated Reading Time:** ~12 minutes
- **Number of Chapters:** 6
- **Transcript Available:** Yes (analyzed from video content)

## Table of Contents

1. Reward Modeling in Reinforcement Learning
  2. Learning from Human Preferences
  3. Key Mathematical Concepts
  4. Visual Elements from the Video
  5. Self-Assessment for This Video
  6. Key Takeaways from This Video
- 

## Video Overview

This video lecture provides a detailed mathematical foundation for **Reward Modeling**, a critical component in modern Reinforcement Learning (RL), particularly for training large-scale Generative AI models. The instructor, Prof. Prathosh A P, explains why an explicit reward model is necessary, the challenges of creating one, and the sophisticated techniques used to learn it from human feedback. The lecture transitions from the theoretical definition of a reward function to the practical challenges of data collection and subjectivity, culminating in a deep dive into the **Bradley-Terry model**, which learns from preferential data rather than absolute scores.

## Learning Objectives

Upon completing this study material, students will be able to: - **Define** the role and mathematical form of a reward function in Reinforcement Learning. - **Understand** why reward modeling is essential for policy gradient algorithms like PPO. - **Identify** the primary challenges in creating reward models, including data acquisition and the subjectivity of human-assigned scores. - **Explain** the concept of **preferential data** as a practical alternative to scalar reward data. - **Describe** the **Bradley-Terry model** for learning a reward function from pairwise comparisons. - **Derive** the loss function used to train a reward model based on the Bradley-Terry framework and log-likelihood maximization.

## Prerequisites

To fully grasp the concepts in this video, students should have a foundational understanding of: - **Reinforcement Learning (RL):** Core concepts such as states ( $S$ ), actions ( $A$ ), policies ( $\pi$ ), and the general purpose of a reward signal. - **Machine Learning:** Basic principles of supervised learning, including regression, loss functions, and optimization using gradient descent. - **Neural Networks:** Familiarity with how neural networks are used as function approximators. - **Basic Probability & Calculus:** Concepts like probability, likelihood, logarithms, and derivatives are used in the derivations.

## Key Concepts Covered

- Reward Function ( $R : S \times A \rightarrow \mathbb{R}$ )

- Supervised Learning for Reward Modeling
  - Preferential Data  $(x, y_w, y_l)$
  - The Bradley-Terry Reward Model
  - Sigmoid Function in Preference Modeling
  - Log-Likelihood Loss Function
- 

## Reward Modeling in Reinforcement Learning

### The Role and Definition of the Reward Function

#### Intuitive Foundation

In Reinforcement Learning, an agent learns to make decisions by interacting with an environment. The central mechanism that guides this learning is the **reward signal**. The reward function is the rulebook that provides this signal.

Imagine training a dog. When the dog performs a desired trick (an “action” in a specific “state”), you give it a treat (a “reward”). If it does something undesirable, you might give it no treat or a negative signal. The reward function in RL is the mathematical formalization of this treat-giving process. It’s a function that looks at the agent’s situation (state) and the action it took, and then assigns a numerical score—the reward—to quantify how “good” that action was. A higher reward encourages the agent to repeat the action in similar situations in the future.

As the instructor highlights (00:17), this reward function is fundamental to almost all policy gradient-based RL algorithms. The agent’s ultimate goal is to learn a policy (a decision-making strategy) that maximizes the total accumulated reward over time. Therefore, having an accurate and well-defined reward function is paramount.

#### Mathematical Definition

(00:39) The reward function, denoted by  $R$ , is formally defined as a mapping from the space of all possible state-action pairs to the set of real numbers.

$$R : S \times A \rightarrow \mathbb{R}$$

Where: -  $S$  is the **state space**, the set of all possible states the agent can be in. -  $A$  is the **action space**, the set of all possible actions the agent can take. -  $S \times A$  is the **Cartesian product** of the state and action spaces, representing the set of all possible state-action pairs. -  $\mathbb{R}$  is the set of **real numbers**, representing the scalar reward value.

In practice, for a given state  $s \in S$  and action  $a \in A$ , the reward function outputs a scalar value  $r(s, a)$ , which tells us the immediate “goodness” of taking action  $a$  in state  $s$  (01:06).

### Training a Reward Model: The Supervised Learning Approach

#### The Ideal (but Impractical) Scenario

(01:20) The instructor explains that if we could easily collect a large dataset of state-action-reward tuples, training a reward model would be a standard supervised learning problem.

- **Dataset:** A collection of tuples  $\{(s_i, a_i, r_i)\}_{i=1}^N$ .
- **Problem:** This can be framed as a **regression task**.
- **Model:** We can use a neural network, parameterized by weights  $\phi$ , to approximate the true reward function. We denote this model as  $r_\phi(s, a)$ .

- **Goal:** Train the neural network  $r_\phi$  such that for any given input pair  $(s, a)$ , its output  $r_\phi(s, a)$  is as close as possible to the true reward  $r$ .
- **Training:** This is typically achieved through **Empirical Risk Minimization (ERM)** (03:19), where we minimize a loss function (e.g., Mean Squared Error) over the dataset using an optimization algorithm like gradient descent.

### Challenges with Scalar Reward Data

(03:27) While the supervised learning approach is simple in theory, it faces significant practical challenges:

1. **Difficulty of Data Annotation:** It is extremely difficult and labor-intensive to obtain a large, high-quality dataset of (**state**, **action**, **reward**) tuples. For complex tasks like generating conversational text, asking a human annotator to assign a precise numerical score (e.g., 7.5 out of 10) to a response is non-intuitive and time-consuming.
2. **Subjectivity and Scale Invariance:** The notion of a “good” response is highly subjective. One annotator might rate a response as an 8/10, while another might rate the same response as a 6/10. This lack of a calibrated and consistent scale across different human annotators makes it hard to create a reliable dataset of scalar rewards.

**Key Insight:** Obtaining absolute, scalar reward scores is difficult. However, humans are much better at making *relative* judgments. It is easier to decide if response A is better than response B than it is to assign absolute scores to both A and B. This insight leads to the use of preferential data.

---

## Learning from Human Preferences

### Preferential Data: A Practical Alternative

(04:43) To overcome the challenges of scalar rewards, modern approaches, especially in training Large Language Models (LLMs), rely on **preferential data**. Instead of asking for a score, annotators are shown an input (e.g., a prompt) and two different model-generated responses. They are simply asked to choose which one they prefer.

This process generates a dataset of triplets:

$$D = \{(x_i, y_{w,i}, y_{l,i})\}_{i=1}^N$$

Where for each data point  $i$ : -  $x_i$  is the input (e.g., a prompt, which corresponds to the state  $s$ ). -  $y_{w,i}$  is the “winning” or **preferred** response (an action). -  $y_{l,i}$  is the “losing” or **non-preferred** response (another action).

This data format is much easier to collect at scale and is more robust to the subjectivity of individual annotators. The following flowchart illustrates the data collection process.

flowchart TD

```

A["Input/Prompt (x)"] --> B["Model Generates Two Responses"];
B --> C["Response 1 (y_a)"];
B --> D["Response 2 (y_b)"];
E["Human Annotator"] --> F["Compares y_a and y_b"];
F -->|y_a is better| G["Record Triplet<br/>(x, y_w=y_a, y_l=y_b)"];
F -->|y_b is better| H["Record Triplet<br/>(x, y_w=y_b, y_l=y_a)"];
G --> I["Preferential Dataset"];
H --> I;

```

**Figure 1:** Flowchart illustrating the collection of preferential data for reward modeling.

## The Bradley-Terry Reward Model

(08:55) Given a dataset of preferences, the next question is how to use it to train a reward model  $r_\phi(x, y)$  that learns to assign a scalar score. The **Bradley-Terry model** provides a classic and powerful framework for this.

### Mathematical Formulation

The core idea is to model the probability that a response  $y_w$  is preferred over a response  $y_l$ . This probability is assumed to be a function of the underlying (and unknown) reward scores of each response.

(10:02) The probability of  $y_w$  being preferred over  $y_l$  for a given input  $x$  is modeled as:

$$P(y_w > y_l | x) = \frac{e^{r_\phi(x, y_w)}}{e^{r_\phi(x, y_w)} + e^{r_\phi(x, y_l)}}$$

- **Intuition:** This formula is a **softmax function** applied to the reward scores of the two responses. It states that the probability of  $y_w$  winning is proportional to its exponentiated reward score relative to the sum of exponentiated scores for both options. If  $r_\phi(x, y_w)$  is much larger than  $r_\phi(x, y_l)$ , this probability will approach 1. If they are equal, the probability is 0.5.

### Simplification using the Sigmoid Function

(11:12) This expression can be elegantly simplified. By dividing the numerator and denominator by  $e^{r_\phi(x, y_w)}$ , we get:

$$P(y_w > y_l | x) = \frac{1}{1 + \frac{e^{r_\phi(x, y_l)}}{e^{r_\phi(x, y_w)}}} = \frac{1}{1 + e^{r_\phi(x, y_l) - r_\phi(x, y_w)}} = \frac{1}{1 + e^{-(r_\phi(x, y_w) - r_\phi(x, y_l))}}$$

This is precisely the definition of the **sigmoid function**,  $\sigma(z) = \frac{1}{1+e^{-z}}$ . Therefore, the probability can be expressed as:

$$P(y_w > y_l | x) = \sigma(r_\phi(x, y_w) - r_\phi(x, y_l))$$

**Key Insight:** The probability of one response being preferred over another is modeled as the sigmoid of the *difference* in their learned reward scores. This elegantly connects the pairwise preference data to the underlying scalar reward model we want to learn.

### The Loss Function for Training

(12:20) To train our reward model  $r_\phi$ , we use the principle of **maximum likelihood estimation**. We want to find the parameters  $\phi$  that maximize the probability of observing the preferences in our dataset  $D$ . This is achieved by maximizing the log-likelihood of the data, which is equivalent to minimizing the negative log-likelihood.

The loss function  $L_{\text{reward}}(\phi)$  is the negative log-likelihood averaged over all triplets in the dataset:

$$L_{\text{reward}}(\phi) = -\mathbb{E}_{(x, y_w, y_l) \sim D} [\log P(y_w > y_l | x)]$$

Substituting our sigmoid expression for the probability, the final loss function becomes:

$$L_{\text{reward}}(\phi) = -\mathbb{E}_{(x, y_w, y_l) \sim D} [\log (\sigma(r_\phi(x, y_w) - r_\phi(x, y_l)))]$$

This loss function is a form of **binary cross-entropy loss**, where the goal is to make the difference  $r_\phi(x, y_w) - r_\phi(x, y_l)$  large and positive.

## Optimization

(14:38) The optimal parameters  $\phi^*$  for the reward model are found by minimizing this loss function with respect to  $\phi$ , typically using stochastic gradient descent (SGD) or its variants (like Adam).

$$\phi^* = \arg \min_{\phi} L_{\text{reward}}(\phi)$$

Once trained, this reward model  $r_{\phi^*}(x, y)$  can provide a scalar reward for any given input and response, which can then be used to optimize a policy model using algorithms like PPO.

---

## Key Mathematical Concepts

- **Reward Function:**

$$R : S \times A \rightarrow \mathbb{R}$$

Assigns a scalar value to a state-action pair, quantifying its immediate desirability.

- **Bradley-Terry Probability Model:**

$$P(y_w > y_l | x) = \frac{e^{r_{\phi}(x, y_w)}}{e^{r_{\phi}(x, y_w)} + e^{r_{\phi}(x, y_l)}}$$

Models the probability of a “winning” response  $y_w$  being preferred over a “losing” response  $y_l$  based on their respective reward scores.

- **Sigmoid Formulation:**

$$P(y_w > y_l | x) = \sigma(r_{\phi}(x, y_w) - r_{\phi}(x, y_l))$$

A more compact and computationally stable representation of the Bradley-Terry model, where  $\sigma(z) = \frac{1}{1+e^{-z}}$ .

- **Reward Model Loss Function:**

$$L_{\text{reward}}(\phi) = -\mathbb{E}_{(x, y_w, y_l) \sim D} [\log(\sigma(r_{\phi}(x, y_w) - r_{\phi}(x, y_l)))]$$

The negative log-likelihood loss used to train the reward model parameters  $\phi$  by maximizing the probability of the observed human preferences.

---

## Visual Elements from the Video

The lecture is presented on a digital whiteboard. The key visual elements are the handwritten mathematical formulas and their explanations.

- **(00:39) Reward Function Definition:** The instructor writes down the formal definition  $R : S \times A \rightarrow \mathbb{R}$  and explains that it takes a state and an action to produce a real number indicating how “good” the action was.
- **(05:41) Preferential Data Format:** The instructor introduces the concept of preferential data, outlining the input  $x$ , the preferred output  $y_w$ , and the non-preferred output  $y_l$ .
- **(08:55) Bradley-Terry Reward Models:** A new section is started to introduce the model for learning from preferential data.
- **(10:02) Bradley-Terry Probability Equation:** The core softmax-like probability equation  $P(y_w > y_l)$  is written and explained.
- **(11:12) Sigmoid Formulation:** The instructor derives the equivalent sigmoid formulation  $\sigma(r_{\phi}(x, y_w) - r_{\phi}(x, y_l))$ , which is a central part of the lecture.

- **(12:20) Reward Loss Function:** The final log-likelihood loss function is formulated, showing how the model is trained by maximizing the likelihood of the observed preferences.
- 

## Self-Assessment for This Video

1. **Question 1:** What is the mathematical definition of a reward function, and what is its primary purpose in Reinforcement Learning?
  2. **Question 2:** Why is it often impractical to train a reward model using a simple supervised regression approach on (state, action, reward) tuples? Name at least two challenges.
  3. **Question 3:** What is “preferential data,” and why is it considered a more practical alternative for collecting human feedback? Describe the format of a single data point.
  4. **Question 4:** Explain the intuition behind the Bradley-Terry model. How does it relate the probability of preferring one response over another to their underlying reward scores?
  5. **Problem 1:** Given the Bradley-Terry probability formula  $P(y_w > y_l | x) = \frac{e^{r_\phi(x, y_w)}}{e^{r_\phi(x, y_w)} + e^{r_\phi(x, y_l)}}$ , show the step-by-step derivation to arrive at the simplified sigmoid form:  $\sigma(r_\phi(x, y_w) - r_\phi(x, y_l))$ .
  6. **Problem 2:** You have a reward model  $r_\phi$ . For a prompt  $x$ , it generates two responses,  $y_1$  and  $y_2$ . The model assigns rewards  $r_\phi(x, y_1) = 2.5$  and  $r_\phi(x, y_2) = 1.0$ . According to the Bradley-Terry model, what is the probability that a human would prefer  $y_1$  over  $y_2$ ?
- 

## Key Takeaways from This Video

- **Reward is Central:** The reward function is the cornerstone of policy gradient RL methods, as it directly or indirectly influences all major components of the learning algorithm.
- **Scalar Rewards are Hard:** Collecting absolute, numerical reward scores from humans is difficult, subjective, and does not scale well.
- **Preferences are Easier:** Humans are much better at making relative judgments (A is better than B) than absolute ones. This makes **preferential data** a more robust and scalable source of human feedback.
- **Bradley-Terry for Preferences:** The Bradley-Terry model provides a principled way to learn a scalar reward function  $r_\phi(x, y)$  from pairwise preference data  $(x, y_w, y_l)$ .
- **Loss as Log-Likelihood:** The reward model is trained by maximizing the log-likelihood of the observed human preferences, which translates to minimizing a binary cross-entropy-style loss function based on the sigmoid of the reward difference.