

# Mathematical Framework for Transition Matrices and Expected Steps in Trait Learning Models with Trait Frequencies

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Model Framework</b>	<b>1</b>
2.1	Directed Graph Representation of Traits . . . . .	1
2.2	Agent Repertoire and Learnability . . . . .	2
2.3	Trait Frequencies . . . . .	2
2.4	Learning Strategies . . . . .	2
2.5	Payoffs . . . . .	2
<b>3</b>	<b>Transition Probabilities</b>	<b>3</b>
3.1	Adjusted Weights . . . . .	3
3.2	Transition Probabilities . . . . .	3
<b>4</b>	<b>Expected Steps Until Absorption and Expected Payoff</b>	<b>3</b>
4.1	Fundamental Matrix . . . . .	3
4.2	Expected Steps . . . . .	3
4.3	Expected Payoff . . . . .	3
4.4	Expected Payoff per Step . . . . .	4
<b>5</b>	<b>Fixed-Point Iteration for Trait Frequencies</b>	<b>4</b>
5.1	Iterative Method . . . . .	4
5.2	Iteration Steps . . . . .	4

## 1 Introduction

This document presents a mathematical framework for modeling trait learning processes in agents using Markov chains, incorporating trait frequencies. The model computes transition matrices, determines the expected number of steps until all traits are learned (absorption), and calculates the expected payoff per step.

## 2 Model Framework

### 2.1 Directed Graph Representation of Traits

The model represents traits and their prerequisite relationships as a directed graph  $G = (V, E)$ :

- $V = \{0, 1, \dots, n - 1\}$  is the set of  $n$  traits.
- $E \subseteq V \times V$  is the set of directed edges, where an edge  $(i, j)$  indicates that trait  $i$  is a prerequisite for trait  $j$ .

The adjacency matrix  $A \in \{0, 1\}^{n \times n}$  encodes the graph structure:

$$A_{ij} = \begin{cases} 1 & \text{if there is a directed edge from trait } i \text{ to trait } j, \\ 0 & \text{otherwise.} \end{cases}$$

The set of parent traits (prerequisites) for trait  $j$  is:

$$P_j = \{i \in V : A_{ij} = 1\}.$$

## 2.2 Agent Repertoire and Learnability

An agent's repertoire is represented by a vector  $r \in \{0, 1\}^n$ , where:

$$r_j = \begin{cases} 1 & \text{if trait } j \text{ is known (learned) by the agent,} \\ 0 & \text{otherwise.} \end{cases}$$

A trait  $j$  is *learnable* at state  $r$  if all its prerequisites are learned:

$$L_j(r, A) = (1 - r_j) \cdot \prod_{i \in P_j} r_i.$$

Thus,  $L_j(r, A) = 1$  if and only if trait  $j$  is unlearned and all its parent traits are learned.

## 2.3 Trait Frequencies

Each trait  $j$  has an associated frequency  $f_j \geq 0$  in the population. The frequencies are initialized and updated through a fixed-point iteration process. The frequency of the first trait (trait 0) is ignored throughout the model.

## 2.4 Learning Strategies

Learning strategies determine the base weights assigned to each trait when deciding which trait to learn next.

**Base Weights** The base weight vector  $\mathbf{w}^* \in \mathbb{R}^n$  is defined as:

$$w_j^* = p_j,$$

where  $p_j \geq 0$  represents the payoff or intrinsic value associated with trait  $j$ .

## 2.5 Payoffs

The payoff for each trait is a function of its distance from the root trait (trait 0) and a parameter  $\alpha$ :

$$p_j = u_j + \alpha d_j,$$

where:

- $u_j \sim U(0, 1)$  is a uniform random variable,
- $d_j$  is the distance of trait  $j$  from the root trait,
- $\alpha \geq 0$  is a parameter that controls the influence of distance on payoff.

### 3 Transition Probabilities

#### 3.1 Adjusted Weights

The adjusted weight assigned to trait  $j$  at state  $r$  is:

$$w_j(r) = \begin{cases} f_j \cdot w_j^*, & \text{if } L_j(r, A) = 1, \\ 0, & \text{otherwise.} \end{cases}$$

#### 3.2 Transition Probabilities

From state  $r$ , the probability of transitioning to state  $r'$  by learning trait  $j$  is:

$$P(r \rightarrow r') = \begin{cases} \frac{w_j(r)}{W(r)}, & \text{if } r' = r + e_j, \\ 0, & \text{if } r' \neq r \text{ and } r' \neq r + e_j \text{ for any } j, \\ 0, & \text{if } W(r) > 0 \text{ and } r' = r, \\ 1, & \text{if } W(r) = 0 \text{ and } r' = r. \end{cases}$$

Here,  $e_j$  is the unit vector with a 1 at position  $j$  and 0 elsewhere, and  $W(r)$  is the total weight of learnable traits at state  $r$ :

$$W(r) = \sum_{k: L_k(r, A)=1} w_k(r).$$

If  $W(r) = 0$  (i.e., there are no learnable traits at state  $r$ ), then:

$$P(r \rightarrow r) = 1,$$

and all other transition probabilities are zero.

### 4 Expected Steps Until Absorption and Expected Payoff

#### 4.1 Fundamental Matrix

Let  $Q$  be the submatrix of the transition matrix  $P$  representing transitions among transient states (states that are not absorbing). The fundamental matrix  $N$  is given by:

$$N = (I - Q)^{-1},$$

where  $I$  is the identity matrix of the same size as  $Q$ .

#### 4.2 Expected Steps

The expected number of steps until absorption starting from the initial state  $r_0$  is:

$$t = \sum_r N_{r_0 r}.$$

#### 4.3 Expected Payoff

Let  $p(r)$  be the expected payoff at state  $r$ . The expected total payoff until absorption is:

$$E[P] = \sum_r N_{r_0 r} \cdot p(r).$$

#### 4.4 Expected Payoff per Step

The expected payoff per step is calculated as:

$$E[P_{step}] = \frac{E[P]}{t}.$$

### 5 Fixed-Point Iteration for Trait Frequencies

#### 5.1 Iterative Method

To determine the trait frequencies  $\{f_j\}$ , a fixed-point iteration method is used. The iteration starts with all frequencies set to 1.0 for  $j = 1, \dots, n-1$ , and the frequency of the first trait is ignored throughout.

#### 5.2 Iteration Steps

1. **Initialization:**

$$f_j^{(0)} = 1.0, \quad \text{for } j = 1, \dots, n-1.$$

2. **Iteration** ( $k \geq 0$ ):

(a) Compute the adjusted weights:

$$w_j^{(k)}(r) = \begin{cases} f_j^{(k)} \cdot w_j^*, & \text{if } L_j(r, A) = 1, \\ 0, & \text{otherwise.} \end{cases}$$

(b) Construct the transition matrix  $P^{(k)}$  using the adjusted weights  $w_j^{(k)}(r)$ .

(c) Compute the fundamental matrix:

$$N^{(k)} = (I - Q^{(k)})^{-1}.$$

(d) Update trait frequencies:

$$f_j^{(k+1)} = \frac{\sum_r N_{r0r}^{(k)} \cdot r_j}{t^{(k)}}, \quad \text{for } j = 1, \dots, n-1,$$

where:

$$t^{(k)} = \sum_r N_{r0r}^{(k)}.$$

3. **Convergence Check:** Repeat until the frequencies converge, i.e.,

$$\max_j |f_j^{(k+1)} - f_j^{(k)}| < \epsilon,$$

for a small tolerance  $\epsilon > 0$ .