

# Optimizing Elevator Control with Reinforcement Learning

DRAFT

**Omar Elbaghdadi**

(2577566)

A thesis presented for the degree of  
Bachelor Econometrics and Operations Research



School of Business and Economics

Vrije Universiteit Amsterdam

The Netherlands

June 25, 2018

# **Thesis Title**

Thesis Subtitle

**Author Name**

## **Abstract**

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

## Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>Reinforcement Learning</b>	<b>5</b>
2.1	Exploration and Exploitation . . . . .	5
2.2	Elements of Reinforcement Learning . . . . .	5
2.3	Markov Decision Process . . . . .	6
2.4	Returns . . . . .	6
2.5	Policies . . . . .	7
2.6	Value Functions . . . . .	7
2.7	Optimality . . . . .	8
<b>3</b>	<b>Model of the System</b>	<b>10</b>
3.1	System Dynamics . . . . .	10
3.2	State Space . . . . .	11
3.3	Action Set . . . . .	12
3.4	Performance Measures and Reward . . . . .	13
3.5	Simulation . . . . .	14
3.6	Traffic and Passenger Arrival . . . . .	14
<b>4</b>	<b>The Algorithm</b>	<b>15</b>
4.1	$Q$ -learning . . . . .	15
4.2	Function Updates . . . . .	16
4.3	Calculating Omniscient Reinforcements . . . . .	16
4.4	Making Decisions and Updating $Q$ -Values . . . . .	16
<b>5</b>	<b>Results and Analysis</b>	<b>18</b>
5.1	Annealing Schedules . . . . .	18
<b>6</b>	<b>Discussion</b>	<b>19</b>
6.1	Simulator . . . . .	19
6.2	Exploration Strategy . . . . .	19
6.3	Action Constraints . . . . .	19
6.4	Multi-Agent . . . . .	20
6.5	Traffic Profiles . . . . .	20
<b>A</b>	<b>Appendix</b>	<b>21</b>
A.1	Boltzmann Instability . . . . .	21
A.2	Solving Return Integral . . . . .	21
A.3	Elevator Motion . . . . .	22

## Summary of Notation

This section outlines the notation used for the rest of the thesis. Capitalized letters are used for random variables throughout.

$:=$	equality relationship that is true by definition
$\Pr(X = x)$	probability of random variable $X$ taking on the value $x$ .
$\Pr(X = x \mid Y = y)$	conditional probability of random variable $X$ taking on the value $x$ given random variable $Y$
$\mathbb{E}(X)$	expected value of random variable $X$
$\max_a f$	largest value of the function $f$ that depends on the variable $a$ .
$\mathbb{R}$	set of real numbers
$\in$	is an element of, e.g. $s \in S$
$\subseteq$	is a subset of, e.g. $X \subseteq S$
$\alpha$	step-size parameter
$\beta$	continuous-time discount-rate parameter
$\gamma$	discrete-time discount-rate parameter
$n$	number of floors in the building considered
$Q$	current estimate of $Q$ -value

## 1 Introduction

(FACTCHECK) Life is hard to imagine without elevators nowadays. According to ??, the number of tall buildings constructed keeps increasing every year. Elevators, which play a vital role as a means of transportation in multi-storied buildings, and their service quality are thus becoming increasingly important. An elevator control system handles passenger calls while trying to optimize service quality, e.g. by minimizing passenger waiting time.

The core of classical and even state-of-the-art elevator control strategies are human-designed heuristics. Human-designed control policies perform sufficiently well most of the time, but they are difficult and costly to design. Not only are they suboptimal, heuristic control strategies are also inflexible, due to their inability to deal with traffic patterns that had been considered when designing the algorithm [11].

To combat these weaknesses, several researchers have proposed the application of learning to elevator control. *Reinforcement Learning* (RL) has been given an especially large amount of research in this [5, 8, 11, 12, 7]. These researchers have shown that the application of RL improves optimality of the control policies compared to traditional heuristic policies. An especially realistic model was tested by [11]. Instead of relying on a pre-designed strategy, RL algorithms learn an optimal control policy at running time. Oftentimes, the elevator system is very complex and hard to model. The difficulty of the task is caused primarily by:

- *state space*: there is an enormous amount of possible combinations of car calls, hall calls, and elevator positions and directions.
- *asynchronicity of events*: hall calls or car calls can be signaled at any time moment
- *nonstationarity*: the rate of incoming calls can change both in the short (in the course of a day) and long (in the course of weeks and months) term.
- *partial observability*: the number of passengers waiting on floors and in the elevator is unknown.

The advantage with RL methods is that it is okay not to know all the dynamics of the environment. It is possible to apply *model-free learning* in that case [10]. It is therefore possible to apply the methods without knowing about any of the traffic patterns in advance. To some extent, it is also possible to keep adapting the policy in the face of changing traffic patterns.

The elevator control problem is interesting for serving as a benchmark to RL algorithms as well, since rewards encoding applicable system performances are easily obtained. Furthermore, we are able to test how well the algorithms behave when dealing with a stochastic environment.

The aim of an elevator controller is to provide adequate service to all passengers in the system while trying to optimize some performance measures (usually minimize the waiting time for passenger calls and the travel time for passenger commands).

## 2 Reinforcement Learning

Reinforcement learning is learning what to do—how to map situations to actions—so as to maximize a numerical reward signal. The learner is not told the correct action to take, but instead must discover which actions yield the most reward by trying them. Actions may affect not only the immediate reward but also the next situation and, through that, all following rewards. The idea that we learn by interacting with our environment is a very intuitive one. Whether we are learning to drive a car or to hold a conversation, we are aware of how our environment responds to our actions, and we want to control what happens through our behavior [10].

A learning problem involves interaction between an active decision-making agent and its environment, in which the agent tries to achieve a goal despite facing uncertainty about its environment. The agent's actions are allowed to affect the future state of the environment, thereby affecting the options and opportunities available to the agent at later times. Correct choice requires taking into account not only the direct effect of the choice, but also indirect, delayed consequences of actions.

All these features are a good fit to modelling an elevator controller task. Note that we use the word controller instead of agent here, but they are interchangeable. The controller is able to sense the state of its environment, which can be described by the elevator's position and which buttons are pressed, among others. The controller is taking sequential actions that affect the position of the elevator, thereby affecting the state of the environment. The controller's goal is to serve passengers as quickly as possible. Its goals relate to the state of the environment.

### 2.1 Exploration and Exploitation

A fundamental challenge that arises in reinforcement learning is the trade-off between exploration and exploitation. To obtain a high reward, an agent must choose actions that it has tried in the past and found to yield a higher reward than others. It must exploit the knowledge gained from past actions. However, to discover such actions, it has to keep exploring and try actions it has not tried before. An appropriate balance between exploration and exploitation is necessary. Usually this is done as follows. At the start of training, exploratory actions are highly encouraged. As training moves along, we gradually decrease the probability of taking an exploratory action.

### 2.2 Elements of Reinforcement Learning

Besides the agent and the environment, there are four important subelements in a reinforcement learning system: a policy, a reward signal, a value function, and, optionally, a model of the environment. Their functions will be addressed in this section.

A model of the environment describes the (probabilistic) behavior of the environment. It allows for predictions about the environment's behavior to be made. For example, given a state and action, the model might predict the next state and next reward.

### 2.3 Markov Decision Process

The concepts of reinforcement learning, continuous interaction between an agent and its environment, are formalized as a Markov decision process (MDP). The agent selects actions using some strategy and the environment reacts by presenting the agent new situations. The environment also emits rewards. The agent seeks to maximize its rewards over time by adapting its strategy. See figure 1.

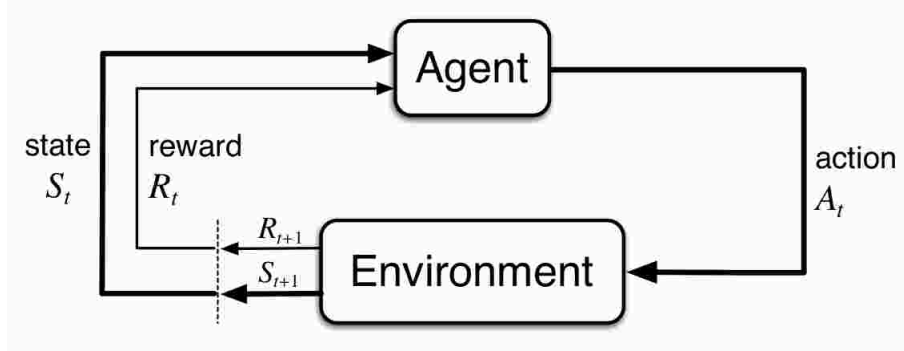


Figure 1: The agent-environment interaction in a Markov decision process [10].

In a discrete-time MDP, the agent and environment interact at each of a sequence of discrete time steps,  $t = 0, 1, 2, 3, \dots$ . At each time step  $t$ , the agent receives some representation of the environment's state,  $S_t \in \mathcal{S}$ , and selects an action,  $A_t \in \mathcal{A}(s)$  depending on the observed state. One time step later, partly as a result of its action, the agent receives a reward,  $R_{t+1} \in \mathcal{R} \subseteq \mathbb{R}$ , and finds itself in a new state,  $S_{t+1}$ .

In a finite MDP, the sets of states, actions, and rewards ( $\mathcal{S}$ ,  $\mathcal{A}$ , and  $\mathcal{R}$ ) all have a finite number of elements. In this case, the random variables  $R_t$  and  $S_t$  have a well defined discrete probability distributions dependent only on the previous state and action. That is, for particular values of these random variables,  $s' \in \mathcal{S}$  and  $r \in \mathcal{R}$ , the probability of those values occurring at time  $t$  given values of the previous state and action is given by

$$p(s', r | s, a) := \Pr(S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a)$$

for all  $s' \in \mathcal{S}$ ,  $r \in \mathcal{R}$ , and  $a \in \mathcal{A}(s)$ .

The probabilities given by the function  $p$  completely determine the dynamics of a finite MDP. From it, anything else we might want to know about the environment, such as the state-transition probabilities  $p(s' | s, a)$ , can be computed.

### 2.4 Returns

A reward signal defines the goal in a reinforcement learning problem. The reward signal thus defines what the good and bad events for the agent are. The agent's objective is to maximize the total reward it receives over the long run. If an action selected by the policy is followed by low reward, the policy may be changed to select some other action in that same situation in the future.

Let  $R_{t+1}, R_{t+2}, R_{t+3}, \dots$  denote the sequence of rewards received after time step  $t$ . Generally, we want to maximize the expected return  $G_t$ , which is some specific function of the reward sequence. For our specific task, we define this to be the sum of discounted future rewards

$$G_t := R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (1)$$

where  $\gamma \in [0, 1]$  is a parameter called the *discount rate*.

The discount rate determines the present value of future rewards: a reward received  $k$  time steps in the future is worth only  $\gamma^{k-1}$  times what it would be worth if it were received immediately. If  $\gamma < 1$ , the infinite sum in (1) has a finite value as long as the reward sequence  $\{R_k\}$  is bounded. As  $\gamma$  approaches 1, the objective takes future rewards into account more strongly.

Returns at successive time steps are related recursively:

$$\begin{aligned} G_t &:= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots \\ &= R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \dots) \\ &= R_{t+1} + \gamma G_{t+1}. \end{aligned}$$

## 2.5 Policies

A policy defines the learning agent's way of behaving, i.e. which actions it selects, at a given time. Roughly speaking, a policy is a mapping from states of the environment to actions to be taken when in those states. The policy alone is sufficient to determine the agent's behavior.

Formally, a policy  $\pi$  is a mapping from states to a probability distribution over the possible actions that can be selected in those states. If the agent is following policy  $\pi$  at time  $t$ , then  $\pi(a | s)$  is the probability that  $A_t = a$  if  $S_t = s$ . Reinforcement learning methods specify how the agent's policy is changed as a result of its experience.

## 2.6 Value Functions

The value of a state  $s$ ,  $v(s)$ , is the return that the agent is expected to acquire in the future, starting from that state. Whereas rewards determine the immediate benefit of an environmental state, values indicate the long-term desirability of states after taking into account the states that are likely to follow, and the rewards available in those states. For example, a state might always yield a low immediate reward but still have a high value because it is often followed by other states that yield high rewards.

The value of a state  $s$  under a policy  $\pi$ , denoted  $v_\pi(s)$ , is the expected return when starting in  $s$  and following  $\pi$  thereafter. For MDPs, we can define  $v_\pi$  formally by

$$v_\pi(s) := \mathbb{E}_\pi[G_t | S_t = s] = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right], \text{ for all } s \in S \quad (2)$$



where  $\mathbb{E}_\pi[\cdot]$  denotes the expected value of a random variable given that the agent follows policy  $\pi$ , and  $t$  is any time step. We call the function  $v_\pi$  the state-value function for policy  $\pi$ . Similarly, we define the value of taking action  $a$  in state  $s$  under a policy  $\pi$ , denoted  $q_\pi(s, a)$ , as the expected return starting from  $s$ , taking the action  $a$ , and thereafter following policy  $\pi$ :

$$q_\pi(s, a) := \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a] = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]. \quad (3)$$

We call  $q_\pi$  the action-value function for policy  $\pi$ . If all the dynamics of an environment were completely specified, exact value functions could be computed. If that were not the case, the value functions  $v_\pi$  and  $q_\pi$  would have to be estimated from (simulated) experience. The first is called a model-based approach, while the latter is called model-free. Since the dynamics of an elevator system are quite complex, we adopt the model-free approach.

Of course, if there are a lot of states, it may not be practical to keep around values for each state individually. The agent would have to parameterize the value functions  $v_\pi$  and  $q_\pi$  (with fewer parameters than states) and adjust the parameters to better match the observed returns. Linear models and artificial neural networks are often used in this case. Similar states would then get similar values. This can also produce accurate estimates, depending on the approximation function. We try to keep the state space small enough to allow for use of tabular methods, since approximate solutions are a bit more complex to analyze.

For any policy  $\pi$  and any state  $s$ , the following condition holds between the value of  $s$  and the value of its possible successor states:

$$v_\pi(s) = \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) \left[ r + \gamma v_\pi(s') \right], \text{ for all } s \in \mathcal{S} \quad (4)$$

which are called the Bellman equations for  $v_\pi$ .

## 2.7 Optimality

Solving a reinforcement learning task means, roughly, finding a policy that achieves a lot of reward over the long run. A policy  $\pi$  is defined to be better than or equal to a policy  $\pi'$  if its expected return is greater than or equal to that of  $\pi'$  for all states. In other words,  $\pi \geq \pi'$  if and only if  $v_\pi(s) \geq v_{\pi'}(s)$  for all  $s \in \mathcal{S}$ . There is always at least one policy that is better than or equal to all other policies. This is an optimal policy. Although there may be more than one, we denote all the optimal policies by  $\pi_*$ . They share the same state-value function, called the optimal state-value function, denoted  $v_*$ , and defined as

$$v_*(s) := \max_{\pi} v_\pi(s), \quad (5)$$

for all  $s \in \mathcal{S}$ . Optimal policies also share the same optimal action-value function, denoted  $q_*$ , and defined as

$$q_*(s, a) := \max_{\pi} q_\pi(s, a), \quad (6)$$

for all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}(s)$ . For the state-action pair  $(s, a)$ , this function gives the expected return for taking action  $a$  in state  $s$  and thereafter following an optimal policy.

The optimal value functions  $v_*(s)$  and  $q_*(s, a)$  can be found by solving the *Bellman optimality equations*. These are given by

$$v_*(s) = \max_a q_{\pi_*}(s, a) = \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')] \quad (7)$$

for  $v_*(s)$ , and

$$q_*(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma \max_a q_*(s', a)] \quad (8)$$

for  $q_*(s, a)$ , for all  $s \in \mathcal{S}$  and for all  $a \in \mathcal{A}(s)$ .

For finite MDPs, a solution to the Bellman optimality equations always exists.

Having  $q_*$  makes choosing optimal actions easy. For any state  $s$ , it can simply find any action that maximizes  $q_*(s, a)$ . It provides the optimal expected long-term return for each state-action pair. The optimal action-value function allows optimal actions to be selected without having to know anything about the environment's dynamics.

Many reinforcement learning methods can be seen as approximately solving the Bellman optimality equations, using actual experienced transitions instead of knowledge of the expected transitions. The on-line nature of reinforcement learning makes it possible to approximate optimal policies in ways that puts more effort into learning to make good decisions for frequently encountered states, at the expense of less effort for infrequently encountered states.

### 3 Model of the System

This section introduces the model of the elevator system studied here. The modelling approach is similar to the one used by [11, 5]. We use this model in all further analysis.

#### 3.1 System Dynamics

The system dynamics are parameterized as follows (as in [9]):

- Number of floors: 5.
- Number of elevator cars: 1.
- Elevator maximum speed: 2.54 m/s.
- Elevator floor time (time it takes for an elevator to pass a floor at full speed): 1.45 seconds.
- Elevator acceleration time (time it takes for an elevator accelerate from a stop to full speed or decelerate from full speed to a stop): 3.595 seconds.
- Elevator load time (the time it takes for a passenger to enter or exit the elevator): 1 second.
- Floor height: 3.66 meters.
- Capacity of the cars: 20 passengers.

Passengers can arrive on any floor, where they can press either the up or down button. We call this a **hall call**. After an elevator arrives at the passenger's floor, they get in and press a button to a floor in the same direction as the elevator. We call this a **car call**.

An elevator system can be very complex, with events such as passenger and elevator arrivals happening asynchronously. Even though time is continuous, it is still possible to model the elevator system as a discrete event system, where significant events happen at discrete times, with the time between them taking on a continuous value. In this case, the constant discount factor  $\gamma$  used in most discrete-time reinforcement learning algorithms is inappropriate. Instead, we take into account the amount of time between events using a time varying discount factor [2]. This is also done to take into account the asynchronicity of events that is well suited to the elevator model. We define returns as the total discounted reward received over time:

$$\int_0^{\infty} e^{-\beta\tau} r_{\tau} d\tau \quad (9)$$

where we integrate over continuous time, instead of the usual discrete-time reinforcement:

$$\sum_{t=0}^{\infty} \gamma^t r_t.$$

Here,  $r_t$  is the immediate cost in the system at discrete time  $t$ ,  $r_{\tau}$  is the instantaneous cost at continuous time  $\tau$  and  $\beta > 0$  is a parameter used to control the rate of exponential decay. A system with a higher  $\beta$  cares more about rewards that are closer in time. The environment is

thusly modeled as a discrete-time Markov decision process with a finite number of states and actions. Using this approach, the states and actions can be fully specified.

The model parameters have been chosen such that the state space is finite and reasonably sized. This is of great importance when using algorithms in which it is necessary to store values for each state, so called tabular methods. If the state space is too large, running time and memory capacity can become intractable. We would need to approximate the value functions.

It is not possible to observe the full state of the system. After a button is pressed, the elevator controller does not know if another passenger arrived after the first one. One way to deal with this is assuming *omniscience*. We assume that in every state, the elevator controller knows how many passengers are waiting at a floor, when they arrived and where they are going. This can then be used to calculate the reinforcement signals. An important thing to note here is that it is not the controller that is receiving this extra information, it is the critic *evaluating* the controller. This information is only important in the training phase. Once the controller is trained, it can be implemented in a real system without needing the extra knowledge.

Another possibility is to let the system learn using only information that would be available to the system on-line. The arrival times of passengers after the first would have to be estimated, so that expected costs can still be computed. We use omniscient reinforcements, since they are slightly more accurate and results will not differ by a large margin [4].

### 3.2 State Space

An elevator controller maps state information to decisions. The definition of the state is therefore of great importance. It needs to properly reflect the current state of the system, measured by quantities that can be practically observed. This includes:

- the state of hall calls and their waiting times,
- awaiting car calls and their waiting times for all elevators,
- position of all elevators,
- moving direction of all elevators,
- velocity of all elevators,
- calls to which particular elevators are allocated.

If we assume the state is defined as above, it is possible to estimate the size of the resulting state space. Let  $n$  be the number of floors. The size of the state space can be estimated as follows:

- at most  $n - 1$  up hall calls:  $2^{n-1}$  possibilities,
- at most  $n - 1$  down hall calls:  $2^{n-1}$  possibilities,
- at most  $n$  car calls for each elevator:  $2^n$  possibilities,
- possible locations for each elevator:  $n$  possibilities,
- possible elevator directions for each elevator: 3 possibilities (up, down or stopped),

which gives a state space size of  $2^{n-1} \cdot 2^{n-1} \cdot 2^n \cdot n \cdot 3$  in a situation where each elevator is considered as a separate learner. The amount of states rises tremendously in the number of floors. The problem already gets computationally intractable for a relatively few number of floors. To combat this and further reduce the size of the state space, we aggregate the states into a more compact representation, while hopefully keeping the most relevant information. Instead of receiving information regarding exactly which hall and car calls are active, the controller only receives:

- the number of up and down awaiting hall calls higher and lower than the elevator's current position,
- the number of car calls to floors in the current moving direction.

Taking into account this particular aggregation of states, we can recalculate the size of the state space as follows:

- the number of remaining **up** hall calls from floors **higher** than the current position: at most  $n - 1$  possibilities,
- the number of remaining **down** hall calls from floors **higher** than the current position: at most  $n - 1$  possibilities,
- the number of remaining **up** hall calls from floors **lower** than the current position: at most  $n - 1$  possibilities,
- the number of remaining **down** hall calls from floors **lower** than the current position: at most  $n - 1$  possibilities,
- the number of remaining car calls in the current moving direction: at most  $n - 1$  possibilities,
- the current position:  $n$  possibilities,
- the current moving direction: 3 possibilities,

which make for a total of  $(n-1)^4 \cdot (n-1) \cdot n \cdot 3$  possible states. This is a vastly smaller amount than before the aggregation. In fact, it is not exponential in  $n$  anymore. Such a size would allow for us to use tabular methods instead of value function approximators, given moderate sizes of  $n$  and  $m$ . Note also that a large subset of states will (almost) never be visited under normal conditions, reducing the effective state space size even further.

### 3.3 Action Set

Now that we have the state space defined, we can move on to defining the actions we can take. What action the elevator is able to take will depend on the state of the system.

We define the actions as follows:

- If the elevator is moving, it can either
  - stop at the next floor, or
  - continue past the next floor.
- If the elevator is not moving, it can either:

- go up, or
- go down.

The system considered is event-based. There are two main types of events. Events of the first type concern waiting times, such as passenger arrivals and transfers in and out of cars. Events of the second type are car arrival events, which are potential decision points for the elevator controllers. A car moving between floors generates a car arrival event when it reaches a point at which it must decide whether to stop at the next floor, or continue past it. Sometimes, the controlling agent is restricted to take a certain action. If a passenger wants to get off at the next floor, it has to stop at the next floor. A potential decision point is only considered a decision point if the choice of action is not constrained.

This approach equals that of [5, 4], but differs from the approach used in [11]. They suggest that instead of using the 2 aforementioned actions, an elevator chooses a floor to serve and commits to that action, under a couple of other constraints. This has wider applicability to practical elevator systems, where the stopping distance for an elevator moving at full speed is longer than half of the distance between floors. Another advantage is the ability to announce the arrival of the car several seconds beforehand. Passengers are given more time to go towards the entrance of the elevator. Although it has practical advantages, it makes the learning task more difficult. Instead of 2 possible actions, at most  $n$  actions can be taken at each state.

In an attempt to build in some basic prior knowledge, we add additional restrictions on what actions can be taken:

- When at the bottom and top floors, we can not choose the action go down and go up respectively. We cannot continue past the bottom and top floors.
- A car cannot turn until it has served all the car calls in its current direction.
- A car can stop at a floor only if there is a call from this floor or there is a car call to this floor.
- Given a choice between moving up and down, it should prefer to move up (since the down-peak traffic tends to push the cars toward the bottom of the building).

### 3.4 Performance Measures and Reward

We need a way to measure the performance of the method we are applying. The goal is to minimize some function of passengers' waiting time. We consider the average passenger waiting time, which is generally considered a primary objective [6]. The waiting time of a passenger is defined as the time between the passenger's arrival at the floor and the passenger's entry into a car. Other possible performance measures are system time and the fraction of passengers waiting more than  $T$  seconds, where  $T$  is typically 60. System time is defined as the waiting time combined with the passenger's travel time.

In a Reinforcement Learning setting, reinforcement is usually formulated as maximizing some reward. In this case, however, it is more convenient to talk about minimizing costs of some sort. We want to define these costs such that minimizing them will lead to a lower average waiting time. The cost function is based on the definition of a reinforcement for an elevator system as in [4, 5].

Since time between events is continuous, we consider instantaneous costs at continuous time moments. For  $\tau \in [t_1, t_2]$ , where  $t_1$  is the time of taking an action for an elevator (i.e. deciding to continue past the next floor or stop) and  $t_2$  is the time when the next decision is required, the reinforcement  $r_\tau$  is defined as follows:

$$r_\tau = \sum_{p \in P} (\tau - t_1 + w_p)^q \quad (10)$$

where  $q$  is any positive real,  $P$  is the set of passengers waiting at time  $t_2$  and  $w_p$  is the amount of time a passenger  $p \in P$  has already waited at time  $t_1$ . It should be noted that  $w_p$  does not depend on  $\tau$ . Special care is needed to handle any passengers that begin or end waiting between  $t_1$  and  $t_2$ , which is addressed in section 4.2

We choose to set  $q = 2$ . A quadratic function has nice properties and makes it easy to compute the return integral. Furthermore, if  $q$  were lower, very high waiting times would not be punished as much.

### 3.5 Simulation

The controller learns to act from experience. This experience is generated by simulation. We simulate an hour of elevator traffic and call that an episode. Every episode, the environment is reset to an initial state and passenger arrivals are generated. We can distinguish between two types of episodes, training and testing. In training episodes, we still take actions with suboptimal  $q$  values to keep exploring and estimating these values. In testing episodes, however, the training is done and only actions that are regarded as optimal will be taken. Further simulation specifics can be found in Appendix (??)

### 3.6 Traffic and Passenger Arrival

It is important to take into account the traffic profile at a time. General building traffic profiles have been identified [6]. Four important profiles are up-peak, down-peak, inter-floor, and lunchtime. We will concern ourselves only with the down-peak traffic profile. Down-peak is a traffic pattern in which passengers are primarily moving down to the ground floor. An example is people going home at the end of a business day in an office building.

We model the arrival of passengers as a Poisson process with rates that fluctuate throughout the day. Table 1 shows the expected number of passengers arriving at each floor in 5-minute intervals. Of these arriving passengers, 90% are headed for the ground floor, while the rest is inter-floor traffic.

Table 1: Arrival rates for the down-peak traffic profile.

Time	00	05	10	15	20	25	30	35	40	45	50	55
Rate	0.25	0.5	1	1	4.5	3	2	1.75	4.5	1.25	0.75	0.5

## 4 The Algorithm

There are various methods and philosophies to solving Reinforcement Learning problems. The basic idea for most algorithms is to estimate reward values for states. These values are then used to make sure we take actions so that we stay in high-value states.  $Q$  is a function that maps state-action pairs  $(s, a)$  to a numerical value. This value represents the expected total discounted return after taking action  $a$  in state  $s$ . These values are then used to approximate the optimal policy  $\pi_*$ .

### 4.1 $Q$ -learning

We use the  $Q$ -learning algorithm [3], which can be categorized as an off-policy Temporal Difference (TD) control method [10], where the  $Q$ -update is generally defined by

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[ r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right] \quad (11)$$

$$= (1 - \alpha)Q(s_t, a_t) + \alpha \left[ r_{t+1} + \gamma \max_a Q(s_{t+1}, a) \right]. \quad (12)$$

It is, in essence, a stochastic approximation to the Bellman Optimality equations. The advantage of this method is that we do not need to fully specify a model of the environment. We do not need to know, for instance, state transition probabilities and the distribution of reinforcements. These can be difficult to specify in an environment as complex as the elevator system. Instead, we use a simulator to generate episodes.

$Q$ -learning is an off-policy method, because it does not necessarily learn the same policy that is being followed. The learned action-value function,  $Q$ , directly approximates the optimal action-value function  $q_*$ , independent of the policy being followed. The policy being followed still determines which state-action pairs visited and updated. However, all that is required for correct convergence is that all pairs continue to be updated. Under this assumption and a variant of the usual stochastic approximation conditions on the sequence of step-size parameters  $\alpha_t$ ,  $Q$  has been shown to converge to  $q_*$  with probability 1.

The  $Q$ -learning algorithm is given below.

---

**Algorithm 1**  $Q$ -learning (off-policy TD control) for estimating  $\pi \approx \pi_*$

---

- 1: Initialize  $Q(s, a)$ , for all  $s \in S$ ,  $a \in A(s)$ , arbitrarily, and  $Q(\text{terminal} - \text{state}, \cdot) = 0$
  - 2: **repeat**
  - 3:   Initialize  $S$
  - 4:   **repeat**
  - 5:     Choose  $A$  from  $S$  using policy derived from  $Q$
  - 6:     Take action  $A$ , observe  $R$ ,  $S'$
  - 7:      $Q(S, A) \leftarrow (1 - \alpha)Q(S, A) + \alpha \left[ R + \gamma \max_a Q(S', a) \right]$
  - 8:      $S \leftarrow S'$
  - 9:   **until** ( $S$  is terminal)
  - 10: **until** episodes end
-



## 4.2 Function Updates

We need to modify Algorithm 1 to meet our requirements. We now minimize over costs instead of maximizing over rewards. Since we are acting in continuous time, it is necessary to change the  $Q$  function update rule described in (12). For a time interval  $[t_1, t_2]$  as considered above, the update rule for the  $Q$ -value for state  $s_{t_1}$  and action  $a_{t_1}$  is given at time  $t_2$  by

$$Q(s_{t_1}, a_{t_1}) \leftarrow (1 - \alpha)Q(s_{t_1}, a_{t_1}) + \alpha \left[ r_{[t_1, t_2]} + e^{-\beta(t_2 - t_1)} \min_a Q(s_{t_2}, a) \right] \quad (13)$$

where the total discounted cost for time interval  $[t_1, t_2]$ ,  $r_{[t_1, t_2]}$ , is given by

$$\int_{t_1}^{t_2} e^{-\beta(\tau - t_1)} \sum_{p \in P} (\tau - t_1 + w_p)^2 d\tau. \quad (14)$$

which can be solved by parts (see Appendix A.1) to yield

$$\sum_{p \in P} \left( \frac{2}{\beta^3} + \frac{2w_p}{\beta^2} + \frac{w_p^2}{\beta} - e^{-\beta(t_2 - t_1)} \left[ \frac{2}{\beta^3} + \frac{2(w_p + t_2 - t_1)}{\beta^2} + \frac{(w_p + t_2 - t_1)^2}{\beta} \right] \right). \quad (15)$$

## 4.3 Calculating Omniscient Reinforcements

We have to be careful in doing the accounting of reinforcements. To deal with discontinuities due to passenger arrivals and transfers, we update the cost between decisions incrementally. In the omniscient reinforcement scheme, which we are using, accumulated cost for the car is updated after every passenger arrival event (when a passenger arrives at a floor), passenger transfer event (when a passenger gets on or off of a car), and when a control decision is made. The car has a cost storage  $R$ , where the total discounted cost it has incurred since its last decision, at time  $d$ , is accumulated. At the time of each aforementioned event, the following computations are performed: Let  $t_0$  be the time of the last event and  $t_1$  the time of the current event. For each passenger  $p$  that has been waiting between  $t_0$  and  $t_1$ , let  $w_0(p)$  and  $w_1(p)$  be the total time that passenger  $p$  has waited at  $t_0$  and  $t_1$  respectively. Then for each car  $i$ ,

$$\Delta R = \sum_{p \in P} e^{-\beta(t_0 - d)} \left[ \frac{2}{\beta^3} + \frac{2w_0(p)}{\beta^2} + \frac{w_0^2(p)}{\beta} \right] - e^{-\beta(t_1 - d)} \left[ \frac{2}{\beta^3} + \frac{2w_1(p)}{\beta^2} + \frac{w_1^2(p)}{\beta} \right]. \quad (16)$$

To keep reinforcements from blowing up, they are scaled down by a factor  $10^6$ .

## 4.4 Making Decisions and Updating Q-Values

A car generates a car arrival event when it is moving between floors and reaches a point at which it must decide to either continue or stop at the next floor. Sometimes, the controller is constrained to a single action, e.g. stopping when a passenger wants to get off at the next floor. In such a case,

the event is not seen as a decision point and the  $Q$ -values will not be updated. The algorithm used by the agent for making decisions and updating its  $Q$ -value estimates is as follows:

1. At time  $t_1$ , observing state  $s_{t_1}$ , the car arrives at a decision point. It selects an action  $a$  according to the Boltzmann distribution over its  $Q$ -value estimates:

$$\Pr(\text{stop}) = \frac{e^{Q(s_{t_1}, \text{continue})/T}}{e^{Q(s_{t_1}, \text{continue})/T} + e^{Q(s_{t_1}, \text{stop})/T}}$$

where  $T$  is a positive “temperature” parameter that is annealed (decreased) during learning.  $T$  controls the amount of exploration in the selection of actions. At the beginning of learning, when the  $Q$ -value estimates are very inaccurate, higher values of  $T$  are used, which give almost equal probabilities to each action. Later in learning, when the  $Q$ -value estimates are more accurate, lower values of  $T$  are used, which give higher probabilities to actions that are thought to be superior, while still allowing some exploration to gather more information about the other actions. It is important to consider how to decrease  $T$  over time, which is discussed in Section 5.1.

2. Let the next decision point for the car be at time  $t_2$  in state  $s_{t_2}$ . After the car has updated its  $R$  value as described above, it adjusts its estimate of  $Q(s_{t_1}, a)$  toward the following target value:

$$R + e^{-\beta(t_2-t_1)} \min_{\{\text{stop}, \text{cont}\}} Q(s_{t_2}, \cdot).$$

3. Let  $s_{t_1} \leftarrow s_{t_2}$  and  $t_1 \leftarrow t_2$ . Go to step 1.

## 5 Results and Analysis

### 5.1 Annealing Schedules

In section 4.2, a strategy for choosing actions was outlined. We sample from the Boltzmann distribution over our  $Q$ -values with temperature parameter  $T$  controlling the randomness over time. As training progresses,  $T$  is annealed according to the following schedule:

$$T = 2 \cdot factor^h$$

where  $factor$  is a positive number less than 1 that determines how quickly the temperature decreases, while  $h$  is the number of episodes simulated until now. The number of episodes to train for depends on  $factor$ . We stop training only when the temperature reaches a predetermined (small) number  $T_{end}$ . Rewriting the above equation in terms of number of hours yields

$$h = \log_{factor} \left( \frac{T}{2} \right).$$

Assuming  $T_{end} = 0.001$  and  $factor = 0.9995$ , we would thus train for  $\log_{0.9995}(0.0005) \approx 15198$  episodes.

Not only do we anneal temperature, we also anneal update step-size parameter  $\alpha$  according to the following schedule:

$$\alpha = 0.01 \cdot 0.99975^h.$$

## 6 Discussion

omniscience vs on-line

### 6.1 Simulator

To train the algorithm, a simulator was built. It is practically impossible to build a simulator that replicates the real world exactly. To deal with motion of the elevator, for instance, many details like how to achieve acting forces on the elevator were left aside. Passenger arrivals are usually modeled while by a Poisson process, although this does not need be the true process. The Poisson assumption simplifies the simulation of arrivals though, as interarrival times can be drawn from an exponential distribution. We also assume that the time it takes for passengers to enter or exit the elevator is 1 second. Other approaches model loading time as a random variable with a mean 1 truncated Erlang distribution. We sample the state of the environment every 0.01 seconds and update accordingly. It can be imagined that sampling more frequently increases accuracy of calculations, but increases simulation running time. The sampling frequency is chosen while weighing these two factors against each other. The specifics of the modelling of elevator motion are described in Appendix A.2.

### 6.2 Exploration Strategy

To make sure we keep exploring actions that seem suboptimal at a current timestep, we use the Boltzmann distribution over the  $Q$ -values (Section 4.3). One of the other commonly used exploration strategies is called  $\epsilon$ -greedy. Using this strategy, an agent takes the (what it thinks to be) optimal action with probability  $1 - \epsilon$  and assigns probability  $\epsilon$  to the other actions. The disadvantage of this method compared to Boltzmann action selection is that it does not take into account the degree to which  $Q$ -values differ.

There are shortcomings to the Boltzmann strategy as well though. It estimates a measure of how optimal the agent thinks the action is, not how certain it is about that optimality. While this is useful information, it is not exactly what would best aid exploration. What we really want to understand is the agent's uncertainty about the value of different actions. There are other strategies such as Bayesian approaches that perform better, but are more complex to implement.

### 6.3 Action Constraints

Focusing the experience of the learning agent onto the most appropriate areas of the state space is important [1]. Constraining the possible actions to take is done exactly for this reason. Here we discuss the choice of some of the action constraints and possible implications and improvements.

- *A car can stop at a floor only if there is a call from this floor or there is a car call to this floor.*

This constraint was imposed when modelling a multi-agent setting. However, in the case of 1 car, it may substantially worsen performance (**CHECK THIS OUT FURTHER**). If the car

was moving down and a new passenger arrived above it while no passengers are under it, it would have to continue on to the ground floor before being able to turn. We look at the implications of removing this constraint.

- *Given a choice between moving up and down, it should prefer to move up.*

This constraint was chosen in the interest of having less actions to learn, which decreases complexity and learning time. Since the traffic profile studies in this case, down-peak traffic, tends to push the cars toward the bottom of the building, we assume it is generally more preferable to move upwards. This changes when we deal with traffic profiles that have more interfloor traffic. In that case, it becomes more important to learn whether moving up or down is the better choice.

## 6.4 Multi-Agent

To keep complexity at bay, we have only considered a building with a single elevator. A building with more elevators corresponds to a multi-agent environment. Multi-agent environments are harder to deal with, because of added stochasticity and non-stationarity due to the changing stochastic policies of the other agents.

## 6.5 Traffic Profiles

We have trained the controller only on the down-peak traffic profile. In a real elevator system, it would have to be trained on other traffic profiles, such as up-peak and interfloor, as well. Other actions would have to be introduced as well. For instance, it would be useful to have actions to open and close the doors in up-peak traffic. It may also be possible to then learn when to switch between traffic profiles.

## A Appendix

### A.1 Boltzmann Instability

We select actions by computing the Boltzmann distribution over  $Q$ -values in a state. However, computing

$$\frac{e^{Q_1/T}}{e^{Q_1/T} + e^{Q_2/T}}$$

is numerically unstable. The fraction will be imprecise, since the value of both the denominator and numerator will be extremely large. At a certain point it becomes impossible to even calculate the exponentials. We use a little trick to remedy this. Let  $Q_* = \max(Q_1/T, Q_2/T)$ . Then we compute

$$\frac{e^{Q_1/T - Q_*}}{e^{Q_1/T - Q_*} + e^{Q_2/T - Q_*}}$$

which is mathematically equivalent to the above but numerically stable.

### A.2 Solving Return Integral

To compute returns, we need to solve the integral

$$\begin{aligned} & \int_{t_1}^{t_2} e^{-\beta(\tau - t_1)} \sum_{p \in P} (\tau - t_1 + w_p)^2 d\tau \\ &= \sum_{p \in P} \int_{t_1}^{t_2} e^{-\beta(\tau - t_1)} (\tau - t_1 + w_p)^2 d\tau. \end{aligned}$$

Looking only at the inside integral and substituting  $\tau_*$  for  $\tau - t_1$ , we get

$$\int_0^{t_2 - t_1} e^{-\beta\tau_*} (\tau_* + w_p)^2 d\tau_*.$$

We first solve the indefinite integral by parts:

$$\begin{aligned} \int e^{-\beta\tau_*} (\tau_* + w_p)^2 d\tau_* &= -\frac{1}{\beta} e^{-\beta\tau_*} (\tau_* + w_p)^2 + 2 \int \frac{1}{\beta} e^{-\beta\tau_*} (\tau_* + w_p) d\tau_* \\ &= -\frac{1}{\beta} e^{-\beta\tau_*} (\tau_* + w_p)^2 + 2 \left( -\frac{1}{\beta^2} e^{-\beta\tau_*} (\tau_* + w_p) + \int \frac{1}{\beta^2} e^{-\beta\tau_*} d\tau_* \right) \\ &= -\frac{1}{\beta} e^{-\beta\tau_*} (\tau_* + w_p)^2 + 2 \left( -\frac{1}{\beta^2} e^{-\beta\tau_*} (\tau_* + w_p) - \frac{1}{\beta^3} e^{-\beta\tau_*} \right). \end{aligned}$$

Evaluating this expression at the endpoints:

$$\begin{aligned} & \left[ -\frac{1}{\beta} e^{-\beta\tau_*} (\tau_* + w_p)^2 + 2 \left( -\frac{1}{\beta^2} e^{-\beta\tau_*} (\tau_* + w_p) - \frac{1}{\beta^3} e^{-\beta\tau_*} \right) \right]_0^{t_2 - t_1} \\ &= \frac{w_p^2}{\beta} + \frac{2w_p}{\beta^2} + \frac{2}{\beta^3} - e^{-\beta(t_2 - t_1)} \left( \frac{(t_2 - t_1 + w_p)^2}{\beta} + \frac{2(t_2 - t_1 + w_p)}{\beta^2} + \frac{2}{\beta^3} \right) \end{aligned}$$

Adding the sum over passengers yields the expression in (15):

$$\sum_{p \in P} \left( \frac{2}{\beta^3} + \frac{2w_p}{\beta^2} + \frac{w_p^2}{\beta} - e^{-\beta(t_2-t_1)} \left[ \frac{2}{\beta^3} + \frac{2(t_2-t_1+w_p)}{\beta^2} + \frac{(t_2-t_1+w_p)^2}{\beta} \right] \right)$$

### A.3 Elevator Motion

Although in a real world system acceleration times and stopping times can depend on the state of the elevator, we constrain ourselves to fixed acceleration times. Acceleration of an elevator is modeled by a sine wave, with velocity and position its first and second antiderivatives respectively. When the elevator chooses to stop while still accelerating, halfway in between floors, its deceleration is modeled by a parabola. Parameters are chosen such that, when an elevator stops at a floor, its velocity and acceleration are both 0. An example motion trajectory is given in Figure 2.

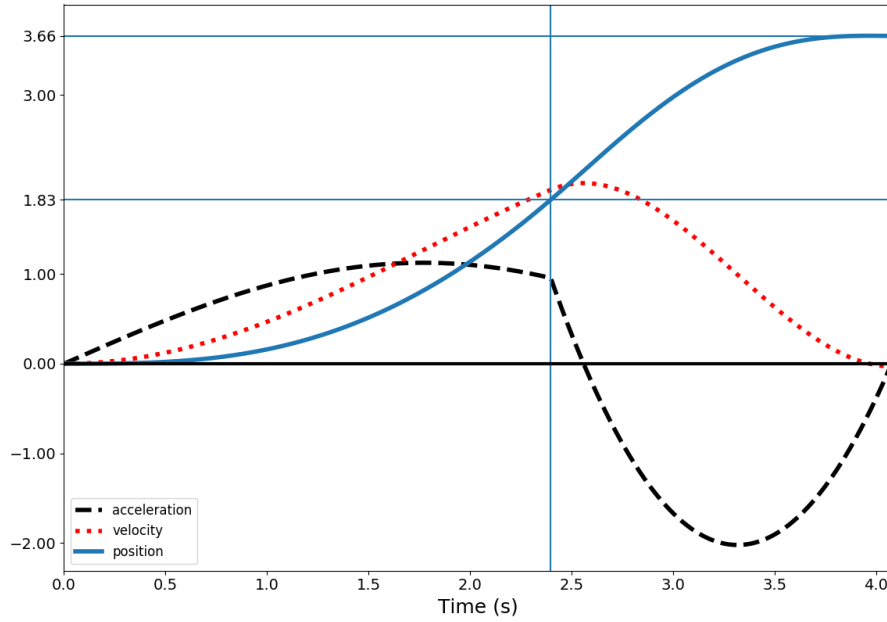


Figure 2: Acceleration, velocity and position trajectories when elevator decides to stop halfway between floors. Note that floor height is 3.66 meters.

## References

- [1] Andrew G. Barto, S. J. Bradtke, and Satinder P. Singh. Learning to act using real-time dynamic programming. *Artificial Intelligence*, 72(1):81–138, 1995. rtdp.
- [2] Steven J. Bradtke and Michael O. Duff. Reinforcement learning methods for continuous-time markov decision problems. In Gerald Tesauro, David S. Touretzky, and Todd K. Leen, editors, *NIPS*, pages 393–400. MIT Press, 1994.
- [3] C. J. C. H. Watkins. *Learning from Delayed Rewards*. PhD thesis, King’s College, Cambridge, 1989.
- [4] Robert H. Crites and Andrew G. Barto. Elevator group control using multiple reinforcement learning agents. *Machine Learning*, 33:235, 1998.
- [5] Robert H. Crites and Andrew G. Barto. Improving elevator performance using reinforcement learning. April 16 1998.
- [6] James Lewis. A dynamic load balancing approach to the control of multiserver polling systems with applications to elevator system dispatching. 1991.
- [7] H. Li. The implementation of reinforcement learning algorithms on the elevator control system. In *2015 IEEE 20th Conference on Emerging Technologies Factory Automation (ETFA)*, pages 1–4, Sept 2015.
- [8] David L. Pepyne and Christos G. Cassandras. Optimal dispatching control for elevator systems during uppeak traffic. *IEEE Trans. Contr. Sys. Techn*, 5(6):629–643, 1997.
- [9] David L. Pepyne, Douglas P. Looze, Christos G. Cassandras, and Theodore E. Djaferis. Application of q-learning to elevator dispatching. *IFAC Proceedings Volumes*, 29(1):4742 – 4747, 1996. 13th World Congress of IFAC, 1996, San Francisco USA, 30 June - 5 July.
- [10] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: an Introduction*. The MIT Press, 2012.
- [11] Tomasz Walczak and Pawel Cichosz. A distributed learning control system for elevator groups. In Leszek Rutkowski, Ryszard Tadeusiewicz, Lotfi A. Zadeh, and Jacek M. Zurada, editors, *Artificial Intelligence and Soft Computing - ICAISC 2006, 8th International Conference, Zakopane, Poland, June 25-29, 2006, Proceedings*, volume 4029 of *Lecture Notes in Computer Science*, pages 1223–1232. Springer, 2006.
- [12] Xu Yuan, Lucian Buşoniu, and Robert Babuška. Reinforcement learning for elevator control. *IFAC Proceedings Volumes*, 41(2):2212 – 2217, 2008. 17th IFAC World Congress.