

---

# EXPLORING THE PARETO FRONT OF MULTI-OBJECTIVE COVID-19 MITIGATION POLICIES USING REINFORCEMENT LEARNING

---

A PREPRINT

**Mathieu Reymond**  
Vrije Universiteit Brussel  
Brussels, Belgium  
mathieu.reymond@vub.be

**Conor F. Hayes**  
National University of Ireland Galway  
Galway, Ireland

**Lander Willem**  
University of Antwerp  
Antwerp, Belgium

**Roxana Rădulescu**  
Vrije Universiteit Brussel  
Brussels, Belgium

**Steven Abrams**  
Hasselt University  
Hasselt, Belgium

**Diederik M. Roijers**  
HU University of Applied Sciences Utrecht  
Utrecht, the Netherlands

**Enda Howley**  
National University of Ireland Galway  
Galway, Ireland

**Patrick Mannion**  
National University of Ireland Galway  
Galway, Ireland

**Niel Hens**  
Hasselt University  
Hasselt, Belgium

**Ann Nowé**  
Vrije Universiteit Brussel  
Brussels, Belgium

**Pieter Libin**  
Vrije Universiteit Brussel  
Brussels, Belgium

## ABSTRACT

Infectious disease outbreaks can have a disruptive impact on public health and societal processes. As decision making in the context of epidemic mitigation is hard, reinforcement learning provides a methodology to automatically learn prevention strategies in combination with complex epidemic models. Current research focuses on optimizing policies with respect to a single objective, such as the pathogen’s attack rate. However, as the mitigation of epidemics involves distinct, and possibly conflicting, criteria (i.e., prevalence, mortality, morbidity, cost), a multi-objective decision approach is warranted to learn balanced policies. To lift this decision-making process to real-world epidemic models, we apply deep multi-objective reinforcement learning and build upon a state-of-the-art algorithm, Pareto Conditioned Networks (PCN), to learn a set of solutions that approximates the Pareto front of the decision problem. We consider the first wave of the Belgian COVID-19 epidemic, which was mitigated by a lockdown, and study different deconfinement strategies, aiming to minimize both COVID-19 cases (i.e., infections and hospitalizations) and the societal burden that is induced by the applied mitigation measures. We contribute a multi-objective Markov decision process that encapsulates the stochastic compartment model that was used to inform policy makers during the COVID-19 epidemic. As these social mitigation measures are implemented in a continuous action space that modulates the contact matrix of the age-structured epidemic model, we extend PCN to this setting. We evaluate the solution set that PCN returns, and observe that it correctly learns to reduce the social burden whenever the hospitalization rates are sufficiently low. In this work, we thus demonstrate that multi-objective reinforcement learning is attainable in complex epidemiological models and provides essential insights to balance complex mitigation policies.

**Keywords** Multi-objective reinforcement learning · Epidemic control · COVID-19 epidemic models

# 1 Introduction

As shown by the COVID-19 pandemic, infectious disease outbreaks represent a major global challenge. To this end, understanding the complex dynamics that underlie these epidemics is essential. Epidemiological transmission models allow us to capture and understand such dynamics and facilitate the study of prevention strategies through simulation. However, developing efficient mitigation strategies remains a challenging process, given the non-linear and complex nature of epidemics.

To address these challenges, reinforcement learning provides a methodology to automatically learn mitigation strategies in combination with complex epidemic models [Libin et al., 2021]. Previous research focused on optimising policies with respect to a single objective, such as the pathogen’s attack rate, while the mitigation of epidemics is a problem that inherently covers distinct and possibly conflicting criteria (i.e., prevalence, mortality, morbidity, cost). Therefore, optimizing on a single objective requires that these distinct criteria are somehow aggregated into a single metric. Generally, during learning, a decision maker will rely on an expert’s assistance to manually design such a metric. However, a decision maker may be unaware of all possible trade-offs between different solutions, which may leave the decision maker unaware of possible solutions that better reflect their preferences, potentially resulting in sub-optimal performance. Furthermore, manually designing such metrics is time consuming, costly and error-prone, as this non-intuitive process requires repetitive and tedious tuning to achieve the desired behaviour [Hayes et al., 2021]. Moreover, taking a single objective approach has several other limiting factors [Hayes et al., 2021, Roijers et al., 2013].

To alleviate the challenging process of defining the learning problem explicitly, we can directly learn optimal behaviours by explicitly taking a multi-objective approach. By assuming that a decision maker will always prefer solutions for which at least one objective improves, it is possible to learn a set of optimal solutions called the *Pareto front*. This enables decision makers to review each solution on the Pareto front before making a decision, thereby being aware of the trade-offs that a solution may imply.

In this work, we investigate the use of *multi-objective reinforcement learning* (MORL) to learn a set of solutions that approximate the Pareto front of multi-objective mitigation strategies. We consider the first wave of the Belgian COVID-19 epidemic, which was mitigated by a hard lockdown [Willem et al., 2021]. When the incidence of confirmed cases were steadily dropping, epidemiological experts were asked to investigate strategies to exit from the stringent lockdown which was imposed. Here, we consider the epidemiological model developed by Abrams et al. [2021] that was constructed to describe the Belgian COVID-19 epidemic, and was fitted to hospitalisation incidence data and serial sero-prevalence data. This model constitutes a stochastic discrete-time age-structured compartmental model that simulates mitigation strategies by varying social distancing parameters concerning school, work and leisure contacts. Based on this model, we contribute MOBElCov, a novel multi-objective epidemiological environment, in the form of a multi-objective Markov decision process (MOMDP) [Roijers et al., 2013]. MOBElCov encapsulates the epidemiological model developed by Abrams et al. [2021] to implement state transitions, with an action space that combines a proportional reduction of school, work and leisure contacts at each time step. Furthermore, it defines a reward function based on two objectives: the attack rate (i.e., proportion of the population affected by the pathogen) and the social burden that is induced by the mitigation measures.

To learn and explore the trade-offs between the attack rate and social burden we present a state-of-the-art MORL approach, Pareto Conditioned Networks (PCN) [Reymond et al., 2022]. PCN uses a single neural network to learn the policies that belong to the Pareto front. As PCN is an algorithm designed for discrete action-spaces, we extend it towards continuous action-spaces to accommodate MOBElCov’s action-space. With this continuous action variant of PCN, we explore the Pareto front of multi-objective COVID-19 mitigation policies.

By evaluating the solution set of mitigation policies returned by PCN, we observe that PCN minimises the social burden in scenarios where hospitalization rates are sufficiently low. Therefore, in this work we illustrate that multi-objective reinforcement learning can be utilised to provide important insights surrounding the trade-offs between complex mitigation policies in real-world epidemiological models.

## 2 Background

### 2.1 Multi-Objective Reinforcement Learning

Many real-world decisions problems often have multiple and possibly conflicting objectives. Multi-objective reinforcement learning (MORL) can be used to find optimal solutions for sequential decision making problems with multiple objectives [Hayes et al., 2021]. For MORL, problems are typically modelled as a multi-objective Markov decision process (MOMDP), i.e., a tuple,  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{T}, \gamma, \mathcal{R})$ , where  $\mathcal{S}$ ,  $\mathcal{A}$  denote the state and action spaces respectively,  $\mathcal{T}: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  denotes a probabilistic transition function,  $\gamma$  is a discount factor determining the importance

of future rewards and  $\mathcal{R}: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}^n$  is an  $n$ -dimensional vector-valued immediate reward function, where  $n$  corresponds to the number of objectives.

While MOMDPs define the actions an agent can take, it does not define which actions are taken in each state. We call this the agent’s *policy*. This policy dictates the agent’s behavior. More formally, a policy  $\pi$  is a probabilistic function  $\mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  determining the probability, for each state  $s \in \mathcal{S}$ , to take an action  $a \in \mathcal{A}$ . We measure the performance of  $\pi$  through the expected sum of discounted rewards it achieves. This is denoted as the  $V$ -value:

$$\mathbf{V}^\pi = \mathbb{E} \left[ \sum_{t=0}^h \gamma^t \mathbf{r}_t \mid \pi, s_0 \right], \quad (1)$$

where  $\mathbf{r}_t = \mathcal{R}(s_t, a_t, s_{t+1})$ .

For single-objective RL, i.e., when  $n = 1$ , the goal is to find the policy  $\pi^*$  that maximizes the  $V$ -value:

$$\pi^* = \arg \max_{\pi} V^\pi. \quad (2)$$

In contrast, in MORL,  $n > 1$  which leads to vectorial returns. In this case, there can be policies for which, without any additional information, it is impossible to know if one is better than the other. For example, we cannot decide which policy between  $\pi_1, \pi_2$  is optimal if they lead to expected returns (also called  $V$ -values)  $\mathbf{V}^{\pi_1} = (0, 1)$ ,  $\mathbf{V}^{\pi_2} = (1, 0)$  respectively. We call these solutions *non-dominated*, i.e., solutions for which it is impossible to improve an objective without hampering another. The set that contains all the non-dominated solutions of the decision problem is called the *Pareto front*  $\mathcal{F}$ . Our goal is to find the set of policies that lead to all the  $V$ -values contained in the Pareto front  $\Pi^* = \{ \pi \mid \mathbf{V}^\pi \in \mathcal{F} \}$ . In general, we call any set of  $V$ -values a *solution set*. When a solution set contains only non-dominated  $V$ -values, it is referred to as a *coverage set*. In the case that no  $\pi$  exists that has a  $\mathbf{V}^\pi$  dominating any of the solutions in a coverage set, then this coverage set is the Pareto front.

## 2.2 Multi-Objective Metrics

Comparing the learned coverage sets produced by different algorithms is a non-trivial task, as one algorithm’s output might dominate the other in some region of the objective-space, but be dominated in another. Intuitively, one would generally prefer the algorithm that covers a wider range of decision maker preferences. In this work we use several metrics to evaluate our algorithm’s performance.

A widely used metric in the literature is called the *hypervolume* [Zitzler et al., 2003]. This metric evaluates the learned coverage set by computing its volume with respect to a fixed specified reference point. The reference point is taken as a lower bound on the achievable returns so that the volumes are always positive. By definition, the solutions contained in the Pareto front dominate all other possible solutions. Thus, no other solution can further increase the volume under the Pareto front. This means that the hypervolume is the highest for the Pareto front. One disadvantage of the hypervolume is that it can be difficult to interpret. For example, when working in high-dimensional objective-spaces, adding or removing a single point can drastically change hypervolume values, especially if the point lies close to an extremum of said space.

To counterpoise these shortcomings, we additionally evaluate our work on a different metric called the  $\varepsilon$ -indicator  $I_\varepsilon$  [Zitzler et al., 2003], which measures how close a coverage set is to the Pareto front  $\mathcal{F}$ . Intuitively,  $I_\varepsilon$  shows that any solution of  $\mathcal{F}$  is *at most*  $\varepsilon$  better with respect to each objective  $o$  than the closest solution of the evaluated coverage set:

$$I_\varepsilon = \inf_{\varepsilon \in \mathbb{R}} \{ \forall \mathbf{V}^\pi \in \mathcal{F}, \exists \mathbf{V}^{\pi'} \in \hat{\Pi} : \|\mathbf{V}^\pi - \mathbf{V}^{\pi'}\|_\infty \leq \varepsilon \}, \quad (3)$$

where  $\|\cdot\|_\infty$ , the L-infinity norm, returns the magnitude of the largest entry of a vector.

The main disadvantage of this metric is that we need the true Pareto front to compute it, which in the case of our MOMDP (see Section 3) is unknown. To still gain insights from our learned policies, we approximate the true Pareto front using the non-dominated policies across all runs.

We note that the  $\varepsilon$ -indicator metric is quite pessimistic, as it measures worst-case performance [Zintgraf et al., 2015], i.e., e.g., it will still report bad performance as long as a single point of the Pareto front is not correctly modeled, even if all the other points are covered. As such, we also use the  $I_{\varepsilon-mean}$  [Reymond et al., 2022] which measures the *average*  $\varepsilon$  distance of the solutions in  $\mathcal{F}$  with respect to the evaluated coverage set.

Figure 1 shows a visual representation of the hypervolume and  $\varepsilon$  metrics in two dimensions.

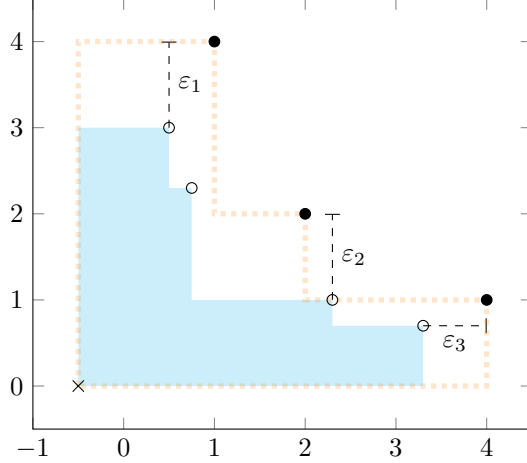


Figure 1: Example of a Pareto front (black dots) and a coverage set (white dots) in a 2-objective environment. The hypervolume metric (in light blue) measures the volume of all dominated solutions with respect to some reference point (cross). The  $\varepsilon$  metrics first compute the maximum distance between each point in the Pareto front and its closest point in the coverage set ( $\varepsilon_i$ ). We can then take their maximum value to compute the  $I_\varepsilon$  metric, or their mean value to obtain the  $I_{\varepsilon-mean}$  metric of the coverage set.

### 3 COVID-19 model and the MOBElCov MOMDP

#### 3.1 Stochastic compartment model for SARS-CoV-2

As an environment to evaluate non-pharmaceutical interventions, we consider the adapted version of the SEIR compartmental model presented by Abrams et al., that was used to investigate exit strategies in Belgium after the first epidemic wave of SARS-CoV-2 [Abrams et al., 2021]. This model constitutes a discrete-time stochastic model, that considers an age-structured population. This model generalises a standard SEIR model<sup>1</sup>, extended to capture the different stages of disease spread and history that are associated with SARS-CoV-2 (i.e., pre-symptomatic, asymptomatic, symptomatic with mild symptoms and symptomatic with severe symptoms) and to represent the stages associated with severe disease, i.e., hospitalisation, admission to the intensive care unit (ICU) and death. As we consider an age-structured population, we consider this extended SEIR structure for  $K = 10$  age groups, i.e.,  $[0 - 10)$ ,  $[10 - 20)$ ,  $[20 - 30)$ ,  $[30 - 40)$ ,  $[40 - 50)$ ,  $[50 - 60)$ ,  $[60 - 70)$ ,  $[70 - 80)$ ,  $[80 - 90)$ ,  $[90, 100+)$ . Contacts of the different age-groups, which impact the propagation rate of the epidemic, are modeled using social contact matrices (SCMs). We define a SCM for 6 different social environments:  $C_{home}$ ,  $C_{work}$ ,  $C_{transport}$ ,  $C_{school}$ ,  $C_{leisure}$ ,  $C_{other}$  for the home, work, transport, school, leisure, other environments respectively. The model is described by a set of ordinary differential equations (ODEs), as described in SI. By formulating this set of differential equations a chain-binomial process (as detailed in SI) we can obtain stochastic trajectories from this model [Abrams et al., 2021].

#### 3.2 Interventions strategies

In order to model different types of interventions, we follow Abrams et al. [2021]. Firstly, in order to consider distinct exit scenarios, we modulate the SCMs to reflect a contact reduction in a particular age group. Secondly, we assume that compliance to the interventions is gradual and model this using a logistic compliance function (see details in SI). We consider a contact reduction function that imposes a proportional reduction of work (including transport)  $p_w$ , school  $p_s$  and leisure  $p_l$  contacts, which is implemented as a linear combination of SCMs:

$$\hat{C} = C_{home} + p_w(C_{work} + C_{transport}) + p_s C_{school} + p_l(C_{leisure} + C_{other}) \quad (4)$$

#### 3.3 The MOBElCov Environment

In order to apply multi-objective reinforcement learning, we construct the MOBElCov MOMDP based on the epidemiological model introduced in Section 3.1 and graphically depicted in Figure 2.

<sup>1</sup>A standard SEIR model divides the population into four different states, i.e., susceptible, exposed, infectious and recovered individuals.

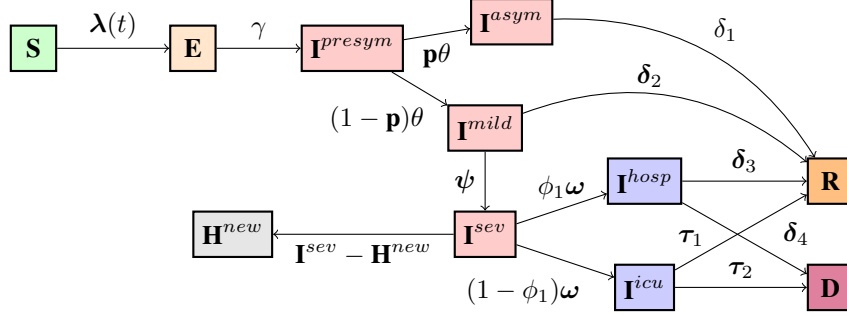


Figure 2: Schematic diagram of the compartmental model for SARS-CoV-2 presented by Abrams et al. [2021] which is used to derive the MOMDP.

**Action-space:** Our actions concern the installment of a social contact matrix  $\hat{C}$ , with a particular reduction (see Section 3.2). To this end, we use the proportional reduction parameters  $p_w, p_s, p_l$  defined in Section 3.2. Thus, each  $\mathbf{a} \in \mathcal{A}$  is a 3-dimensional continuous vector in  $[0, 1]^3$  (i.e.,  $\mathbf{a} = [p_w, p_s, p_l]$ ) which impacts the SCM according to Equation 4.

**Transition function:** The model defined by Abrams et al. [2021] utilises a model transition probability  $M(\mathbf{s}'_m | \mathbf{s}_m, \hat{C})$  (see SI for details on  $M$ ), where  $\hat{C}$  the currently installed SCM, that progresses the epidemiological model in one timestep. We use this function as the transition function in MOBELCov.

In a classical MDP, executing an action  $a_t$  in any state  $s_t$  leads to a next state  $s_{t+1}$  according to the transition function  $\mathcal{T}$ . At every timestep  $t$ , the policy is free to choose the action to perform. In our case, this potentially results in a different restriction  $[p_w, p_s, p_l]$  every week. However, we argue that in the context of mitigation policies, consistency is key, and policies that impose changes too frequently will be hard to follow.

In order to learn realistic and consistent mitigation policies, we incorporate a *budget* on the number of times a policy can change its actions until the terminal state of the MOMDP. Concretely, when the action changes, i.e., if the social restriction proposed by the policy is different from the current one in place, we reduce the budget by one. We only allow action changes as long as there is budget left. We note that we can simulate a no-limit budget setting by setting the budget to the length of the episode.

Finally, for each timestep  $t$ , our transition function  $\mathcal{T}$  uses the model transition probability  $M$  to simulate the model for one week, using  $\hat{C}$  obtained from  $\mathbf{a}_t$ .

**State-space:** The state of the MOMDP concerns a 3-tuple. The first element,  $\mathbf{s}_m$ , directly corresponds to the aggregation of the state variables in the epidemiological model, i.e., a tuple,

$$\{S_k, E_k, I_k^{presym}, I_k^{asym}, I_k^{mild}, I_k^{seu}, I_k^{hosp}, I_k^{icu}, H_k^{new}, D_k, R_k\}, \quad (5)$$

for each age group  $k \in \{1, \dots, K\}$ , where  $S$  encodes the members of the population who are susceptible to infection and  $E$  encodes the members of the population who have been exposed to COVID-19. Moreover,  $I^{presym}, I^{asym}, I^{mild}, I^{hosp}, I^{icu}$  are the members of the population infected with COVID-19 and are, respectively, presymptomatic, asymptomatic, have mild symptoms, are hospitalised, or are in the ICU. Finally,  $H_k^{new}$  represents the number of newly hospitalised individuals in age group  $k$ .

We parameterize the epidemiological model using the mean of the posteriors as specified by Abrams et al. [2021] (details in SI).

The second element of the tuple consists of the social contact matrix  $\hat{C}$  that is currently in place. The reason to incorporate it in the state-space is two-fold. First, Abrams et al. [2021] define a compliance function, simulating the time people need to get used to the new rules set in place. As such, during the simulated week, there is a gradual shift from the current  $\hat{C}$  to the next one. We thus require to know the current  $\hat{C}$ .

Secondly, we require the current  $\hat{C}$  to determine whether the action changes the social restrictions in place – and thus consume part of the budget – or not.

The third element of the tuple concerns of the budget  $\mathbf{b}$ . We note that we incorporate a distinct budget per action-dimension, so  $p_w, p_s$  and  $p_l$  each have their own budget, rendering  $\mathbf{b}$  a vector  $\mathbf{b} = [b_w, b_s, b_l]$ . As such, it is possible that, at timestep  $t$ , the budget for one of the proportional reductions is reduced but not the others.

Therefore, we define a state in MOBElCov as follows:

$$\mathbf{s} = \mathbf{s}_m \cup \hat{C} \cup b \quad (6)$$

**Reward function:** We define a vectorial reward function which considers multiple objectives: attack rate (i.e., infections, hospitalisations) and the social burden imposed by the interventions on the population.

The attack rate in terms of infections is defined as the difference in susceptibles from the current state and next state [Libin et al., 2021]. Since this is a cost that needs to be minimized, we defined the corresponding reward function as the negative attack rate:

$$\mathcal{R}_{\text{ARI}}(\mathbf{s}, \mathbf{a}, \mathbf{s}') = -\left(\sum_{k=1}^K S_k(\mathbf{s}) - \sum_{k=1}^K S_k(\mathbf{s}')\right). \quad (7)$$

The reward function to reduce the attack rate in terms of hospitalisations is defined as the negative number of new hospitalizations:

$$\mathcal{R}_{\text{ARH}}(\mathbf{s}, \mathbf{a}, \mathbf{s}') = -\sum_{k=1}^K H_k^{\text{new}}(\mathbf{s}). \quad (8)$$

Finally, we use the missed contacts resulting from the intervention measures as a proxy for societal burden. To quantify missed contacts, we consider the original social contact matrix  $C$  and the installed social contact matrix  $\hat{C}$  to compute the difference  $\hat{C} - C$ . The resulting difference matrix quantifies the average frequency of contacts missed. To determine missed contacts for the entire population, we apply the difference matrix to the population sizes of the respective age groups that are currently uninfected (i.e., susceptible and recovered individuals). Therefore we define the social burden reward function  $\mathcal{R}_{\text{SB}}$ , as follows:

$$\mathcal{R}_{\text{SB}}(\mathbf{s}, \mathbf{a}, \mathbf{s}') = \sum_{i=1}^K \sum_{j=1}^K (\hat{C} - C)_{ij} S_i(\mathbf{s}) S_j(\mathbf{s}) + \sum_{i=1}^K \sum_{j=1}^K (\hat{C} - C)_{ij} R_i(\mathbf{s}) R_j(\mathbf{s}), \quad (9)$$

where  $S_k(\mathbf{s})$  represents the number of susceptible individuals in age group  $k$  in state  $\mathbf{s}$  and  $R_k$  represents the number of recovered individuals in age group  $k$  in state  $\mathbf{s}$ . In Section 5, we optimize PCN on two different variants for the multi-objective reward function:  $[\mathcal{R}_{\text{ARH}}, \mathcal{R}_{\text{SB}}]$  and  $[\mathcal{R}_{\text{ARI}}, \mathcal{R}_{\text{SB}}]$ .

## 4 Pareto Conditioned Networks

In multi-objective optimization, the set of optimal policies can grow exponentially in the number of objectives. Thus, recovering them all is an expensive process and requires an exhaustive exploration of the complete state space. To address this problem, we use Pareto Conditioned Networks (PCN), a method that uses a single neural network to encompass all non-dominated policies [Reymond et al., 2022]. The key idea behind PCN is to use supervised learning techniques to improve the policy instead of resorting to temporal-difference learning.

**TODO: TODO REPHRASE, FROM PCN PAPER** Using neural networks as function approximators in RL comes with many challenges. One of them is that the target (e.g., the optimal action of the policy) is not known in advance — as opposed to classical supervised learning where the ground-truth target is provided. As the behavior of the agent improves over time, the action used as target can change, often leading to hard-to-tune and brittle learners [??].

Instead of trying to continuously improve the policy by learning actions that should lead to the highest cumulative reward, PCN flips the problem, by learning actions that should lead to any *desired* cumulative reward (be it high or low). In this way, all past trajectories can be reused for supervision, since their returns are known, as well as the actions needed to reach said returns. We can thus train a policy that, conditioned on a desired return, provides the optimal action to reach said return. By leveraging the generalization properties of neural networks, we can accumulate incrementally better experience by conditioning on increasingly higher reward-goals

PCN uses a single neural network that takes a tuple  $\langle \mathbf{s}, \hat{h}, \hat{\mathbf{R}} \rangle$  as input.  $\hat{\mathbf{R}}$  represents the *desired return* of the decision maker, i.e. the return PCN should reach at the end of the episode.  $\hat{h}$  denotes the *desired horizon* that expresses the number of timesteps that should be executed before reaching  $\hat{\mathbf{R}}$ . At execution time, both  $\hat{h}$  and  $\hat{\mathbf{R}}$  are chosen by the decision maker at the start of the episode. Then, at every timestep, the desired horizon is updated according to the perceived reward  $\mathbf{r}_t$ ,  $\hat{\mathbf{R}} \leftarrow \hat{\mathbf{R}} - \mathbf{r}_t$  and the desired horizon is decreased by one,  $\hat{h} \leftarrow \hat{h} - 1$ .

PCN’s neural network has an output for each action  $a_i \in \mathcal{A}$ . Each output represents the confidence the network has that, by taking the corresponding action in  $s$ ,  $\hat{\mathbf{R}}$  will be reached in  $\hat{h}$  timesteps. We can draw an analogy with a classification problem where the network should learn to classify  $(s, \hat{h}, \hat{\mathbf{R}})$  to its corresponding label  $a_i$ .

Similar to classification, PCN requires a labeled dataset with training examples to learn a mapping from input to label. However, contrary to classification, the data present in the dataset is not fixed. PCN collects data from the trajectories experienced while exploring the environment. Thus, the dataset improves over time, as we collect increasingly better trajectories.

Concretely, the PCN algorithm is a loop that repeatedly executes 3 main steps:

1. Use the current policy to interact with the environment and generate trajectories (see Section 4.2).
2. Update the dataset to incorporate these new trajectories and prune the least interesting trajectories in terms of returns (see Section 4.1).
3. Update the policy by training it on the updated dataset (see Section 4.3).

#### 4.1 Building and updating the dataset

Given a trajectory composed of  $T$  transitions  $(s_0, a_0, \mathbf{r}_0, s_1), \dots, (s_{T-1}, a_{T-1}, \mathbf{r}_{T-1}, s_T)$  we can compute, for each transition at timestep  $t$ , the future discounted return  $\mathbf{R}_t = \sum_{i=t}^T \gamma^i \mathbf{r}_i$  and the leftover horizon  $h_t = T - t$ . Since for this trajectory executing action  $a_t$  in state  $s_t$  resulted in return  $\mathbf{R}_t$  in  $h_t$  timesteps, we add a datapoint with input  $\langle s, \hat{h}, \hat{\mathbf{R}} \rangle = \langle s_t, h_t, \mathbf{R}_t \rangle$  and output  $a = a_t$  to the dataset. In other words, when the observed return corresponds to the desired return in that state, then  $a_t$  is the optimal action to take.

Since we have a fix-sized dataset, adding this trajectory means we need to remove another one. Our aim is to keep trajectories that span different parts of the objective space, such that we can explore as many types of trade-offs as possible. However, to improve the trade-offs we currently have in our dataset, we need to only keep solutions that are close to our current estimate of the Pareto front. Thus, we assign a priority for each trajectory in our dataset that combine two metrics: the *crowding distance*, which measures the distance of a solution with its closest neighbors, and the *l2-norm* with the closest non-dominated solution, which indicates how close the solution is from our current estimate of the Pareto front.

The dataset is then updated by only keeping the trajectories with the top- $N$  priorities.

#### 4.2 Generating trajectories with the current policy

In the previous section, we explain how PCN updates the dataset given a trajectory. In this section, we explain how PCN produces said trajectory.

It is unrealistic to expect PCN to reliably produce trajectories with high-valued desired returns when it has only been trained on datapoints originating from random trajectories. Rather, we can expect PCN to produce trajectories with returns in the range of the ones from the current training data. Therefore, if we obtain trajectories reaching high returns, PCN will be able to confidently return high-return policies.

PCN leverages the fact that, due to the generalization capabilities of neural networks, the policies obtained from the network will still be reliable even if the desired return is marginally higher than what is present in the training data. In fact, they will perform similar actions to those in the training data, but lead to a higher return. Thus, we incrementally condition the network on better and better returns, in order to obtain trajectories that extend the boundaries of PCN’s current coverage set.

More precisely, we randomly select a non-dominated return  $\mathbf{R}_{nd}$  and its corresponding horizon  $\hat{h}$  from the dataset. By randomly picking a non-dominated return from the entire coverage set we ensure equal chance of improvement to each part of the objective space. To generate trajectories that improve upon  $\mathbf{R}_{nd}$ , we randomly select an objective  $o$  and increase  $R_{nd,o}$  by a sample from the uniform distribution  $U(0, \sigma_o)$ , where  $\sigma_o$  is the standard deviation for the selected objective, using all non-dominated returns from the trajectories in the dataset. This then becomes our desired return  $\hat{\mathbf{R}}$ . By restricting the improvement to at most  $\sigma_o$ ,  $\hat{\mathbf{R}}$  stays in the range of possible achievable returns and, by only modifying one objective at a time, the changes to the network’s input compared to the training data are kept at a minimum.

### 4.3 Training the network for continuous actions

PCN trains the network as a classification problem, where each class represents a different action. Transitions  $x = \langle \mathbf{s}_t, h_t, \mathbf{R}_t \rangle$ ,  $y = a_t$  are sampled from the dataset, and the ground-truth output  $y$  is compared with the predicted output  $\hat{y} = \pi(\mathbf{s}_t, h_t, \mathbf{R}_t)$ . The predictor (i.e., the policy) is then updated using the cross-entropy loss function:

$$H = - \sum_{a \in \mathcal{A}} y_a \log \pi(a | \mathbf{s}_t, h_t, \mathbf{R}_t) \quad (10)$$

where  $y_a = 1$  if  $a = a_t$  and  $y_a = 0$  otherwise.

While the original PCN algorithm is designed for MOMDPs with discrete actions-spaces, the problem we tackle (see Section 3) is defined in terms of a continuous action-space. We thus extend PCN for the continuous action-space setting. First, we change the output of the neural network such that there is a single output value for each dimension of the action-space. Since the actions should be bound in the domain of possible actions ( $[0, 1]$  in the case of MOBElCov, see Section 3), we apply a tanh non-linearity function on this output. Second, the problem becomes a regression problem instead of a classification problem, as the labeled dataset uses continuous labels  $y = \mathbf{a}_t$  instead of categories. We thus use a Mean Squared Error (MSE) loss to update the policy:

$$MSE = \frac{1}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} (\hat{y}_a - y_a)^2 \quad (11)$$

Since learning the full set of Pareto-efficient policies  $\Pi^*$  requires that the policies  $\pi^* \in \Pi^*$  are deterministic stationary policies [Rojers et al., 2013], we use the output  $\hat{\mathbf{y}}$  as action at execution time. However, PCN improves its policy through exploration, by continuously updating its dataset with better trajectories. Thus, at training time, we follow [Lillicrap et al., 2015] and use a stochastic policy by adding random noise sampled from a Normal distribution to the action:

$$\mathbf{a}_t = \pi(\mathbf{s}_t, h_t, \mathbf{R}_t) + \eta s \text{ with } s \sim \mathcal{N}, \quad (12)$$

where  $\eta$  is a hyper-parameter defining the magnitude of noise to be added.

### 4.4 Coping with stochastic transitions

PCN trains its policy on a dataset that is collected by executing trajectories. It assumes that reenacting a transition from the dataset leads to the same episodic return. When the transition function  $\mathcal{T}$  of the MOMDP is deterministic, the whole trajectory can be faithfully reenacted, which guarantees that we obtain the same return. Combined with the fact that PCN’s policy is deterministic at execution time, conditioning the policy on a target episodic return is equivalent to conditioning it on the V-value  $\mathbf{V}$ .

However, when  $\mathcal{T}$  is stochastic we lose this guarantee. We cope with this issue by adding small random noise to  $\mathbb{R}_t$  when performing gradient descent, which reduces the risk of overfitting [Zur et al., 2009]. Moreover, although the Binomial model of MOBElCov is stochastic, the possible next-states resulting from a state-action pair are similar to each other. This allows PCN to compensate if  $\mathbf{r}_t = \mathcal{R}(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1})$  is somewhat worse than expected. This is confirmed by our experiments (see Section 5), where the coverage sets learned by PCN on the ODE model and on the Binomial model are similar.

## 5 Analysing COVID-19 deconfinement policies

Our goal is to use PCN to learn deconfinement strategies in the MOBElCov environment. We aim to learn policies that balance between the epidemiological objective of minimising the attack rate and the social burden experienced by the population. To this end we consider two cases for the vectorial reward functions  $[\mathcal{R}_{\text{ARH}}, \mathcal{R}_{\text{SB}}]$  and  $[\mathcal{R}_{\text{ARI}}, \mathcal{R}_{\text{SB}}]$ , to learn and analyse policies under different targets with respect to the considered attack rate.

To conduct this analysis, we evaluate our extension of PCN for continuous action-spaces on both the Binomial model. **TODO: reason for binomial.** As explained in Section 4.4, we extend PCN for environments with stochastic transitions. While all results in this Section concern the Binomial model, we observe similar performance on the ODE model, showing that our extension does resolve the stochasticity issue. These results are available in the SI.

Conform to Abrams et al. [2021], the simulation starts on the 1st of March 2020, by seeding a number of infections in the population. Two weeks later, on the 14th of March, the Belgian government initiated a full lockdown of the country. This is implemented by fixing the actions  $p_w, p_s, p_l$  to 0.2, 0, 0.1 respectively. This lockdown ended on the 4th



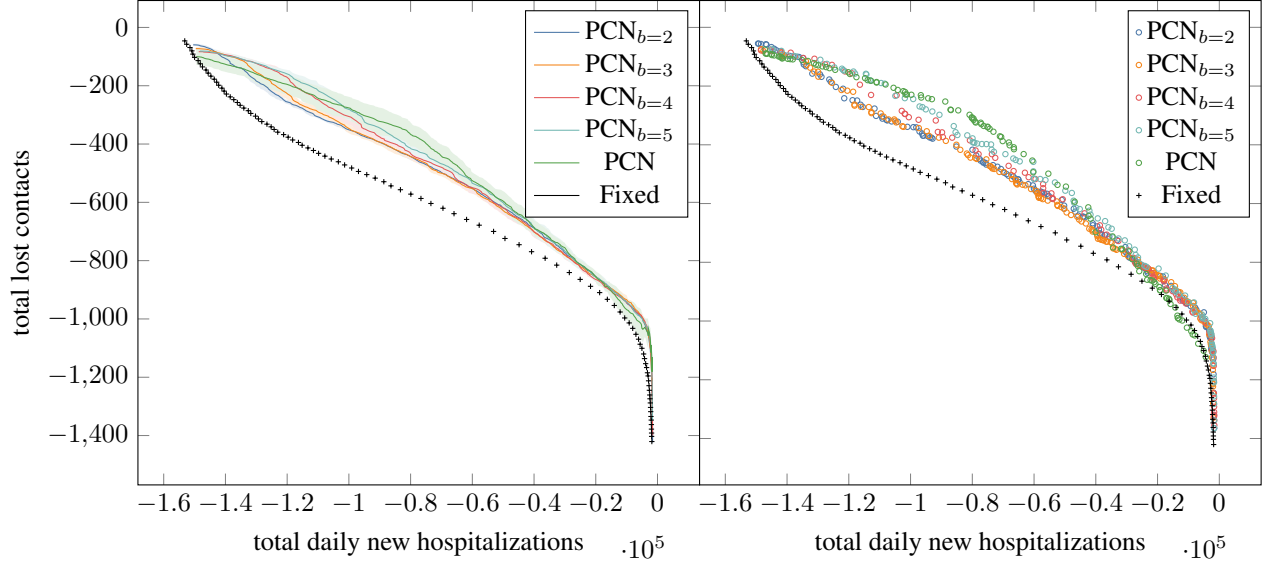


Figure 3: The Pareto front of policies discovered by PCN using the Binomial model, showing the different compromises between the number of hospitalizations and the number of lost contacts. As the budget increases, so does the coverage set learnt by PCN. Changes are most noticeable in the less restrictive trade-offs in terms of social burden.

of May 2020, at which point the government decided on a multi-phase exit strategy to incrementally reduce teleworking, reopen schools and allow leisure activities, such as the re-opening of bars and the cultural sector. It is from this day onward that PCN aims to learn policies for diverse exit strategies, compromising between the total number of daily new hospitalizations and the total number of contacts lost as a proxy for social burden. The simulation lasts throughout the school holidays (from 01/07/2020 to 31/08/2020). Schools are closed during the school holidays, which is simulated by setting  $p_s = 0$ , regardless of the corresponding value outputted by the policy, i.e., during periods of school closure  $p_s$  is ignored.

We draw the analogy with the multi-phase exit strategy established by the Belgian government and the restriction on the number of action-changes imposed by the MOBElCov’s budget  $b$ . Indeed, on the 11th of May 2020, exactly one week after the end of the lockdown, stores and certain types of jobs were allowed to reopen, under strict conditions. This corresponds to altering  $p_w$  in our MOMDP. One week later, on 18th of May, primary and secondary schools reopened for limited sized class-groups, and the cultural sector reopened partially. This is equivalent to increasing  $p_s$  and  $p_l$ . Further changes of restrictions occurred on the 8th of June, 1st, 9th and 25th of July. Thus, we argue that, with a limited budget, we can achieve realistic and implementable policies. In our experiments, we consider budgets of 2 to 5, as these closely relate to the number of changes that occurred until the end of the summer holidays of 2020. As an upper bound, we also consider a no-limit budget setting.

To evaluate the quality of the policies discovered by PCN, we compare it to a baseline. The baseline consists of a set of 100 fixed policies, that iterate over all the possible social restriction levels, with values ranging from 0 to 1. Concretely, each policy uses a fixed proportional reduction  $p_w = p_l = p_s = u$ ,  $u \sim \mathcal{U}(0, 1)$  throughout the episode. In other words, the fixed policies directly operate in a fine-grained manner on the whole contact reduction function  $\hat{C}$ . This allows us to obtain a strong baseline for potential exit strategies over the objective space. We note that while such fixed policies are a feasible approach, they do not scale well in terms of action and objective spaces and they will not be able to provide an adaptive restriction level, which is what we aim to provide using PCN.

All experiments are averaged over 5 runs. The hyper-parameters and the neural network architecture can be found in the SI.

### 5.1 Learned coverage set

We learn a coverage set that ranges from imposing minimal restrictions to enforcing many restrictions (see Figure 4). The coverage set is shown in Figure 3. On the right, we display the coverage set of the best-performing run in terms of

	$[\mathcal{R}_{\text{ARH}}, \mathcal{R}_{\text{SB}}]$			$[\mathcal{R}_{\text{ARI}}, \mathcal{R}_{\text{SB}}]$		
	Hypervolume	$I_\epsilon$	$I_{\epsilon-\text{mean}}$	Hypervolume	$I_\epsilon$	$I_{\epsilon-\text{mean}}$
PCN <sub>b=2</sub>	158.370 $\pm$ 0.811	0.080 $\pm$ 0.011	0.033 $\pm$ 0.002	157.152 $\pm$ 1.023	0.087 $\pm$ 0.006	0.035 $\pm$ 0.002
PCN <sub>b=3</sub>	158.721 $\pm$ 1.439	0.080 $\pm$ 0.012	0.032 $\pm$ 0.003	158.002 $\pm$ 2.081	0.084 $\pm$ 0.009	0.034 $\pm$ 0.005
PCN <sub>b=4</sub>	160.642 $\pm$ 1.582	0.075 $\pm$ 0.007	0.028 $\pm$ 0.003	159.315 $\pm$ 2.601	0.088 $\pm$ 0.018	0.031 $\pm$ 0.006
PCN <sub>b=5</sub>	163.104 $\pm$ 2.386	0.070 $\pm$ 0.023	0.022 $\pm$ 0.005	161.792 $\pm$ 2.464	0.075 $\pm$ 0.015	0.026 $\pm$ 0.005
Fixed	140.479 $\pm$ 0.000	0.139 $\pm$ 0.000	0.073 $\pm$ 0.000	140.479 $\pm$ 0.000	0.139 $\pm$ 0.000	0.073 $\pm$ 0.000
PCN	159.462 $\pm$ 7.713	0.264 $\pm$ 0.115	0.036 $\pm$ 0.020	159.852 $\pm$ 2.395	0.171 $\pm$ 0.093	0.032 $\pm$ 0.006

Table 1: Evaluation metrics for the coverage sets comparing hospitalizations with social burden. In general, an increase of budget results in a better coverage set. Training on infections (ARI) still provides a competitive coverage set in terms of hospitalizations. All PCN coverage sets outperform the baseline.

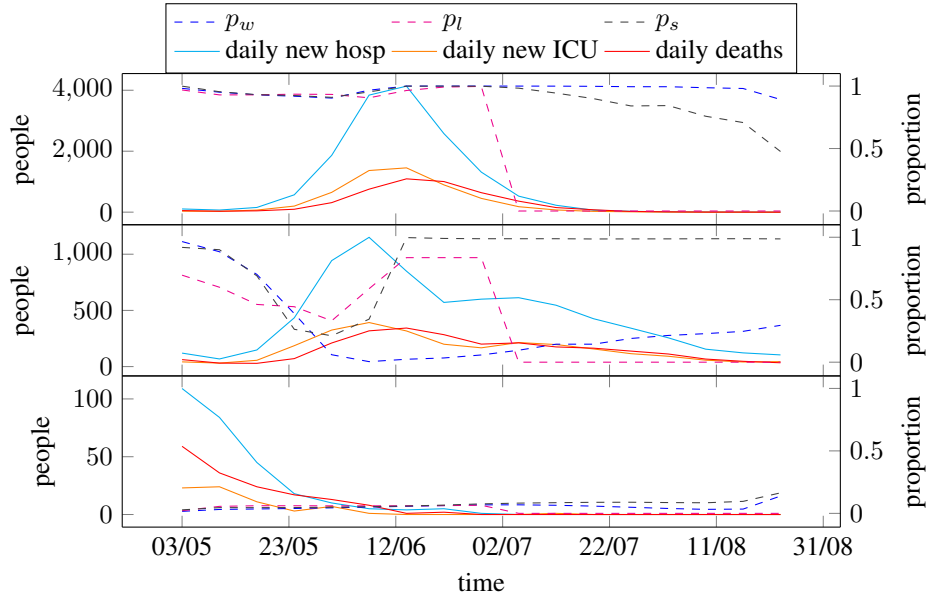


Figure 4: Selection of policies learned by PCN<sub>b=5</sub>, from most restrictive in terms of social burden (top) to least restrictive (bottom).

hypervolume, for each budget setting. On the left, we show an interpolated average of the coverage sets learned by the different runs.

Regardless of the imposed budget, we notice that the coverage sets discovered by PCN almost completely dominate the coverage set of the baseline, demonstrating that there are better alternatives to the fixed policies. This is most visible in the compromising policies, where one has to carefully choose when to remove social restrictions while at the same time minimizing the impact on daily new hospitalizations. In these scenarios, PCN discovers policies that can drastically reduce the total number of new hospitalizations (e.g. more than 20000) for the same social burden. Analysing the executions of such policies (see Figure 4, middle plot) shows a flattened hospitalization curve, with a gradual increase of social freedom during the school holidays such that the peak stays stable and gradually decreases over time.

Interestingly, we notice that the most restrictive policy (i.e. the one that prioritizes hospitalizations over social burden, see Figure 4, bottom plot) still starts to gradually increase  $p_w$  and  $p_l$  from the end of July onward. This is because by then, the epidemic has mostly faded out, and it is safe to reduce social restrictions. The timing of this reduction is important as reducing restrictions too soon can lead to a new wave. PCN learns the impact of its decisions over time, and correctly infers the timing at which restrictions can be safely lifted.

## 5.2 Impact of budgets on the coverage set

Looking at the coverage sets displayed in Figure 3, we notice that the budget impacts the coverage set that has been learnt. In general, an increase of budget results in an increasingly better coverage set. This is to be expected, as a higher budget give the agent more freedom to change its actions as the epidemic progresses. However, we also notice that this increase of coverage sets is not necessarily gradual. Between the run with a budget of 2 and budget of 3, the difference in coverage set seems marginal. **TODO: insight.**

Moreover, we observe that the improvement concentrates around the less restrictive policies in terms of social burden. We believe that this regio contains the most complex policies, as these try to maintain as much social freedom as possible, while containing the number of hospitalisations. We believe that, in these cases, the timing of the actions greatly correlate with the timing and duration of the peak of the epidemic, and a higher budget allows for more fine-grained control about managing this timing. To confirm this, we select, for each budget setting, the solution where the difference in performance is most noticeable, corresponding to the solutions with a total number of hospitalisations around 80000. We plot the execution of the corresponding policies in Figure 5 and analyse the difference they make in terms of social burden. First, we notice that the lower-budget policies are unable to lessen the social restrictions past the peak of the epidemic. In contrast, the setting with no budget restrictions finely controls the restrictions as the epidemic progresses, completely removing restrictions by the end of the wave. We note how the policy with a budget of 5 ressembles the execution of the one without restrictions but that, due to its budget, is unable to progressively lessen the restrictions and instead resorts to a halfway compromise. **TODO: corresponding SB numbers**

Compared to this specific regio of the coverage set, the difference in performance between different budget settings seems marginal around the extremas. **TODO: explain**

Finally, we observe that, while the extremas deliver similar trade-offs regardless of the chosen budget, the results differ for the setting without budget restrictions. Indeed, in this setting, PCN does not seem to learn the most and least restrictive policies, even though there are no constraints on the action-set. As explained in Section 4.2, PCN searches for increasingly better solutions using a stochastic policy. Thus, at every timestep, the action can change compared to the previous one. Continuously outputting the same action (e.g., no social restrictions) becomes a complicated task. In comparison, for the settings with a limited budget, the action stays the same as the previous one once the budget has been spent. It is thus easier to discover the most extreme policies. Thus, in the specific case where we impose no budget, we believe the freedom of action actually hinders PCN’s search.

## 5.3 $\mathcal{R}_{\text{ARH}}$ versus on $\mathcal{R}_{\text{ARI}}$

We now assess the difference in coverage sets when optimizing on  $\mathcal{R}_{\text{ARH}}$  versus  $\mathcal{R}_{\text{ARI}}$ . Although these reward functions have a different scale (there are more infected persons than hospitalised ones), our experiments show that infections and hospitalizations are correlated. This is expected, as during the initial phase of the epidemic, limited immunity was present in the population (few natural immunity and no vaccines), which induces a tight coupling between infection and hospitalisation cases. This is confirmed in Table 1. In this table, we show the different performance metrics (hypervolume,  $I_\varepsilon$ ,  $I_{\varepsilon-\text{mean}}$ ) with respect to the objectives  $[\mathcal{R}_{\text{ARH}}, \mathcal{R}_{\text{SB}}]$ . However, the left-side of the table shows the experiments where these objectives are also used as optimization criterions, while the right-side shows the experiments using  $[\mathcal{R}_{\text{ARI}}, \mathcal{R}_{\text{SB}}]$  as optimization criterions.

Even with the  $\mathcal{R}_{\text{ARI}}$ , the increased budget shows an increase in hypervolume in terms of hospitalisations. Moreover, those hypervolumes are close to the ones trained on  $\mathcal{R}_{\text{ARH}}$ , indicating that their coverage sets are similar. However, regardless of the imposed budget, the hypervolumes are slightly worse. This is to be expected, since those experiments are not directly optimising on  $\mathcal{R}_{\text{ARH}}$ . We make similar conclusions for  $I_\varepsilon$ . For budgets 2, 3 and 5, the difference between the worst-performing policy for the  $\mathcal{R}_{\text{ARI}}$  variant and the  $\mathcal{R}_{\text{ARH}}$  is less than 0.01, indicating less than 1% difference in return values between the two variants when comparing their worst-performing policy. As an exception, we notice that PCN without budget restrictions results in better performance across every metric for the  $\mathcal{R}_{\text{ARI}}$  variant. Still, due to the high standard deviation of the no-budget,  $\mathcal{R}_{\text{ARH}}$  setting, we do not believe this difference is meaningful. **TODO: some kind of meaningful conclusion**

## 5.4 Robustness of policy executions

The dataset of trajectories that PCN is trained on is pruned over time to keep only the most relevant trajectories. The returns of these trajectories are used in Figure 3 to show the discovered coverage set. Each of these returns can be used as desired return for policy execution. We now assess the robustness of the executed policies, by comparing the return obtained after executing the policy with the corresponding target return. For each run, we execute the policy 10 times

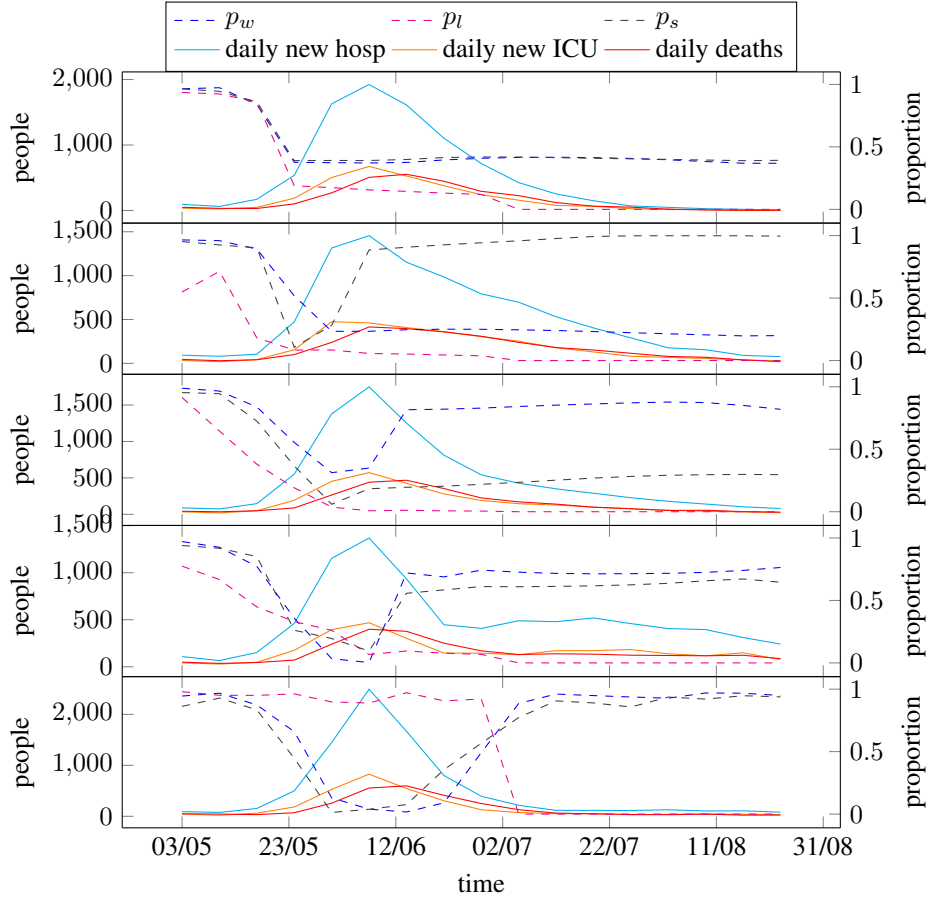


Figure 5: Execution of the policies attaining a number of hospitalisations around 80000, for different budgets. From top to bottom are displayed the policy executions of runs with budget 2,3,4,5 and no-limit, respectively. We notice that the lower-budget policies are unable to lessen the social restrictions past the peak. The setting with no budget restrictions finely controls the restrictions as the epidemic progresses, completely removing restrictions by the end of the wave. Finally, there is no consensus on which social environment to restrict most, some policies providing similar restrictions for  $p_w$  and  $p_l$ , while others putting harsher restrictions on one social environment than the other.

	$I_\epsilon$	$I_{\epsilon-mean}$
$PCN_{b=2}$	$0.047 \pm 0.020$	$0.009 \pm 0.004$
$PCN_{b=3}$	$0.048 \pm 0.022$	$0.007 \pm 0.002$
$PCN_{b=4}$	$0.064 \pm 0.018$	$0.011 \pm 0.003$
$PCN_{b=5}$	$0.058 \pm 0.011$	$0.013 \pm 0.003$
PCN	$0.035 \pm 0.011$	$0.008 \pm 0.004$
Global average	$0.050 \pm 0.017$	$0.010 \pm 0.004$

Table 2: Comparing the difference in the desired return provided to PCN and the actual return PCN obtained when executing its policy. We see that, regardless of the setting, the learned policy faithfully receives a return similar to its desired return.

and compute the  $I_\varepsilon$  and  $I_{\varepsilon-mean}$  metrics with respect to the coverage set learned during the run. We show that the executed policies reliably obtain returns that are similar to the desired return used to condition PCN.

Results are shown in Table 2. The  $I_\varepsilon$  indicators shows that, regardless of the budget, the decision maker will lose at worst a 0.050 normalized return in any of the objectives. On average, it will lose 0.010 normalized returns. Moreover, we emphasize that the learned coverage set contains the non-dominated returns across the whole training procedure. Since the Binomial model is stochasticity, for multiple executions of the same policy, the executions that remains in the coverage set is the one with advantageous transitions due to this stochasticity. Thus, we expect our policy-executions to obtain slightly worse results than the target-solution selected from the coverage set. In light of this, we conclude that the policies trained by PCN are reliable and produce returns as close as possible from their chosen target.

## 6 Related work

Reinforcement learning (RL) has been used in conjunction with epidemiological models to learn policies to limit the spread of diseases and predict the effects of possible mitigation strategies [Probert et al., 2019, Ernst et al., 2006]. For example, RL has been used extensively in modelling and controlling the spread of influenza [Das et al., 2008, Libin et al., 2021, 2018].

RL and Deep RL have been used extensively as a decision making aid to reduce the spread of COVID-19. For example, to learn effective mitigation strategies [Ohi et al., 2020], to learn efficacy of lockdown and travel restrictions [Kwak et al., 2021] and to limit the influx of asymptomatic travellers [Bastani et al., 2021].

Multi-objective methods have also been deployed to learn optimal strategies to mitigate the spread of COVID-19. Wan et al. [Wan et al., 2021] implement a model-based multi-objective policy search method and demonstrate their method on COVID-19 data from China. Given the method is model-based, a model of the transition function must be learned by sampling from the environment. The method proposed by Wan et al. [Wan et al., 2021] only considers a discrete action space which limits the application of their algorithm. Wan et al. [Wan et al., 2021] use linear weights to compute a set of Pareto optimal policies. However, methods which use linear weights can only learn policies on the convex-hull of the Pareto front [Vamplew et al., 2008], therefore the full Pareto front cannot be learned. It is important to note the method proposed by Kompella et al. [Kompella et al., 2020] considers multiple objectives. However, the objectives are combined using a weighted sum with hand-tuned weights which are determined by the authors. The weighted sum is applied by the reward function and a single objective RL method is used to learn a single optimal policy. In contrast to previous work, our approach makes no assumptions regarding the scalarisation function of the user and is able to discover Pareto fronts of arbitrary shape.

## 7 Conclusion and discussion

Making decisions on how to mitigate epidemics has important ethical implications with respect to public health and societal burden. In this regard, it is crucial to approach this decision making from a balanced perspective, to which end we argue that multi-objective decision making is crucial. In this work, we establish a novel approach, i.e., an expert system, to study multi-faceted policies, and this approach shows great potential with respect to future epidemic control. We are aware of the ethical implications that expert systems have on the decision process and we make the disclaimer that all results based on, or derived from, the expert system that we propose should be carefully interpreted by experts in the field of public health, and in a much broader context of economics, well-being and education. We note that the work in this manuscript was conducted by a multi-disciplinary consortium that includes computer scientists and scientists with a background public health, epidemiology and bio-statistics.

In this work, we focus on the clinical outcomes of the intervention strategies and use the reduced contacts as proxy for the quality of life lost. This could be extended into more formal economic evaluations by assessing the health and monetary benefits and costs of different interventions for various stakeholders. The COVID-19 pandemic demonstrated the broad impact of infectious diseases on sectors outside health care. This stresses the need to capture a societal and thus multi-objective perspective in the decision making process on public health and health care interventions. Our learned policies confirm this, showing that focusing solely on reducing the number of hospitalizations results in taking drastic measures – more than a thousand social interactions lost per person over the span of 4 months – that may have a long-lasting impact on the population.

To conclude, we show that multi-objective reinforcement learning can be used to learn a wide set of high-quality policies on real-world problems, providing the decision maker with insightful and diverse alternatives, and showing the impact of taking extreme measures.

## References

- Pieter J. K. Libin, Arno Moonens, Timothy Verstraeten, Fabian Perez-Sanjines, Niel Hens, Philippe Lemey, and Ann Nowé. Deep reinforcement learning for large-scale epidemic control. In Yuxiao Dong, Georgiana Ifrim, Dunja Mladenović, Craig Saunders, and Sofie Van Hoecke, editors, *ECML*, pages 155–170, Cham, 2021. Springer International Publishing. ISBN 978-3-030-67670-4.
- Conor F. Hayes, Roxana Rădulescu, Eugenio Bargiacchi, Johan Källström, Matthew Macfarlane, Mathieu Reymond, Timothy Verstraeten, Luisa M. Zintgraf, Richard Dazeley, Fredrik Heintz, Enda Howley, Athirai A. Irissappane, Patrick Mannion, Ann Nowé, Gabriel Ramos, Marcello Restelli, Peter Vamplew, and Diederik M. Roijers. A practical guide to multi-objective rl and planning, 2021.
- Diederik M Roijers, Peter Vamplew, Shimon Whiteson, and Richard Dazeley. A survey of multi-objective sequential decision-making. *JAIR*, 48:67–113, 2013.
- Lander Willem, Steven Abrams, Pieter JK Libin, Pietro Coletti, Elise Kuylen, Oana Petrof, Signe Møgelmoose, James Wambua, Sereina A Herzog, Christel Faes, et al. The impact of contact tracing and household bubbles on deconfinement strategies for covid-19. *Nature communications*, 12(1):1–9, 2021.
- Steven Abrams, James Wambua, Eva Santermans, Lander Willem, Elise Kuylen, Pietro Coletti, Pieter Libin, Christel Faes, Oana Petrof, Sereina A Herzog, et al. Modelling the early phase of the belgian covid-19 epidemic using a stochastic compartmental model and studying its implied future trajectories. *Epidemics*, 35:100449, 2021.
- Mathieu Reymond, Bargiacchi Eugenio, and Ann Nowé. Pareto conditioned networks. In *Proceedings of the 21st International Conference on AAMAS (2022)*, 2022.
- Eckart Zitzler, Lothar Thiele, Marco Laumanns, Carlos M Fonseca, and Viviane Grunert Da Fonseca. Performance assessment of multiobjective optimizers: An analysis and review. *IEEE Transactions on evolutionary computation*, 7(2):117–132, 2003.
- Luisa M Zintgraf, Timon V Kanters, Diederik M Roijers, Frans Oliehoek, and Philipp Beau. Quality assessment of morl algorithms: A utility-based approach. In *Benelearn 2015: proceedings of the 24th annual ML conference of Belgium and the Netherlands*, 2015.
- Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- Richard M Zur, Yulei Jiang, Lorenzo L Pesce, and Karen Drukker. Noise injection for training artificial neural networks: A comparison with weight decay and early stopping. *Medical physics*, 36(10):4810–4818, 2009.
- William JM Probert, Sandya Lakkur, Christopher J Fonnesbeck, Katriona Shea, Michael C Runge, Michael J Tildesley, and Matthew J Ferrari. Context matters: using reinforcement learning to develop human-readable, state-dependent outbreak response policies. *Philosophical Transactions of the Royal Society B*, 374(1776):20180277, 2019.
- Damien Ernst, Guy-Bart Stan, Jorge Goncalves, and Louis Wehenkel. Clinical data based optimal sti strategies for hiv: a reinforcement learning approach. In *Proceedings of the 45th IEEE Conference on Decision and Control*, pages 667–672. IEEE, 2006.
- Tapas K Das, Alex A Savachkin, and Yiliang Zhu. A large-scale simulation model of pandemic influenza outbreaks for development of dynamic mitigation strategies. *Iie Transactions*, 40(9):893–905, 2008.
- Pieter JK Libin, Timothy Verstraeten, Diederik M Roijers, Jelena Grujic, Kristof Theys, Philippe Lemey, and Ann Nowé. Bayesian best-arm identification for selecting influenza mitigation strategies. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 456–471. Springer, Cham, 2018.
- Abu Quwsar Ohi, MF Mridha, Muhammad Mostafa Monowar, Md Hamid, et al. Exploring optimal control of epidemic spread using rl. *Scientific reports*, 10(1):1–19, 2020.
- Gloria Hyunjung Kwak, Lowell Ling, and Pan Hui. Deep reinforcement learning approaches for global public health strategies for covid-19 pandemic. *PLOS ONE*, 16(5):1–15, 05 2021.
- Hamsa Bastani, Kimon Drakopoulos, Vishal Gupta, Ioannis Vlachogiannis, Christos Hadjicristodoulou, Pagona Lagiou, Gkikas Magiorkinis, Dimitrios Paraskevis, and Sotirios Tsiodras. Efficient and targeted covid-19 border testing via rl. *Nature*, 599(7883):108–113, 2021.
- Runzhe Wan, Xinyu Zhang, and Rui Song. Multi-objective model-based reinforcement learning for infectious disease control. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 1634–1644, 2021.
- Peter Vamplew, John Yearwood, Richard Dazeley, and Adam Berry. On the limitations of scalarisation for multi-objective reinforcement learning of pareto fronts. In *Australasian joint conference on artificial intelligence*, pages 372–378. Springer, 2008.

---

Varun Kompella, Roberto Capobianco, Stacy Jong, Jonathan Browne, Spencer Fox, Lauren Meyers, Peter Wurman, and Peter Stone. Reinforcement learning for optimization of covid-19 mitigation policies. *arXiv preprint arXiv:2010.10560*, 2020.