IMPROVE THE ADAPTATION PROCESS BY REASONING FROM FAILED AND SUCCESSFUL CASES

Anonymous authors

Paper under double-blind review

Abstract

Usually, existing works on adaptation in reasoning-based systems assume that the case base holds only successful cases, i.e., cases having solutions believed to be appropriate for the corresponding problems. However, in practice, the case base could hold failed cases, resulting from an earlier adaptation process but discarded by the revision process. Not considering failed cases would be missing an interesting opportunity to learn more knowledge for improving the adaptation process. This paper proposes a novel approach to the adaptation process in the case-based reasoning paradigm, based on an improved barycentric approach by considering the failed cases. The experiment performed on real data demonstrates the benefit of the method considering the failed cases in the adaptation process compared to the classical ones that ignore them, thus, improving the performance of the case-based reasoning system.

1 INTRODUCTION

Case-based reasoning (CBR) is certainly the most intuitive approach of artificial intelligence to solve a problem since it mimics human behavior in problem-solving. A CBR system looks in its memory represented by a base of previously solved experiments called source cases, for cases having similar problems to the target problem to be solved by adapting their solutions if necessary. The target solution is revised to make sure of its adequacy to solve the target problem and finally the base of cases is enriched following the new experiment of resolution of the target case. Each step of the reasoning process is supported by a process of acquiring the necessary knowledge to perform this step. It is worth highlighting the close connection between the knowledge of the different stages of the CBR approach.

Of the four principal stages of the reasoning process, adaptation is a crucial stage since the quality of the solution heavily depends on its performance. Its focus is on fitting the solutions of similar source cases to meet the specific requirements of the target problem. This is particularly important since the source problems usually do not match the target problem, and as a consequence, without this step, the CBR system cannot ultimately generate an appropriate solution to the target problem.

Awareness of the pivotal role that adaptation plays was noted from the early days of CBR systems, as a result, there is a large number of studies exploring various approaches to acquiring adaptation knowledge to improve its performance. According to Wilke et al. (1997), one can distinguish two approaches to adaptation knowledge acquisition: knowledge-light approaches, which do not require prior knowledge acquisition but exploit the knowledge already contained in the system Petrovic et al. (2016); McDonnell & Cunningham (2006), and knowledge-intensive approaches which need knowledge external to the system, such as knowledge acquired from an expert/user Cordier et al. (2008); Díaz-Agudo & González-Calero (2000); Govedarova et al. (2008).

Existing adaptation approaches focus exclusively on cases whose solutions are deemed relevant to the corresponding problems (hereafter these cases are referred to as successful cases and are denoted by C+). The appreciation of success is subjective to the application domain, e.g., in the context of the CBR application in the elaboration of an energy management system in a building, a successful case would correspond to a scenario satisfying the user's comfort while minimizing the energy expenditure. However, there are also failed cases. A failed case (noted hereafter C-) is a case having an unsatisfactory solution to the problem to solve, in particular, these are cases proposed by the adaptation process but rejected during the validation phase. Moreover, the adaptation process

often involves the acquisition of the knowledge required to generate the adaptation rules. Usually, such knowledge is strongly dependent on the application domain, making the acquisition process complex and challenging to understand and grasp.

Surprisingly, despite a large number of research studies and an increased interest in the adaptation issue, few works are concerned with the challenge of proposing a domain-independent adaptation approach. Even less studies consider adaptation from the solution quality perspective, i.e., addressing both failed and successful cases. These cases are seldom used by the CBR systems even though they constitute potentially useful source of knowledge. We present a comprehensive review of relevant literature on adaptation process in the Appendix A.

In this work, we propose a novel perspective on the adaptation process of the CBR paradigm, based on a fully domain-independent approach and drawing on both successful and failed cases. We argue that the improvement of the performance of the CBR framework relies on the enhancement of the quality of the knowledge supporting the different phases of the CBR approach. In particular, the present study proposes a new approach to the acquisition of adaptation knowledge exploiting both successful cases and failed ones. The approach takes its inspiration from studies in the planning of the path of a robot moving towards a destination in an unknown and insecure environment (includes obstacles). The originality of this approach consists in applying artificial forces to the solution to be proposed to move away from failed source solutions and move closer to successful source solutions.

The rest of this paper is arranged as follows. Section 2 introduces an illustration of motivation and the background of this work. Section 3 details the contribution to harnessing failed and successful cases for a new adaptation approach. An evaluation of the proposed approach is presented and discussed in Section 4, before drawing conclusions about this work and outlining some guidelines for future work in Section 5.

2 MOTIVATING EXAMPLE AND PRELIMINARIES

A CBR-based energy management system (EMS) in a building is a representative case study of the systems relevant to the scope of this study. The objective of an EMS is to fulfill the user's desire for thermal comfort, air quality, etc. while minimizing the energy consumption in the building. Indeed, a building is a complex system whose potential to save energy depends on several factors with dependencies difficult to identify Boulmaiz et al. (2021b), such as climate, building materials, geographical position, and energy rate, but also the occupant of the building exercises a major influence. Findings of earlier works Zhang et al. (2021); Minor & Marx (2017) have already highlighted the advantage of acquiring adaptation knowledge in improving the performance of a CBR-based EMS. Furthermore, due to the growing awareness of environmental issues, several studies have focused on the correlation between energy consumption in a building and the comfort of its occupants, leading to the definition of standards Group (2017); ASHRAE (1992; 2009) to estimate the comfort of users. Thanks to the norms defined in these standards, the revision process can gauge the quality of the target solution proposed by the adaptation process, allowing the retention process to label this solution as a successful case C+ or a failed one C-.

In the CBR-based EMS proposed in Boulmaiz et al. (2021a), the objective is to make the user conscious of the influence of his actions on the energy behavior of the building. For this, the system guides the user in his actions by advising him on a set of actions aiming at decreasing the energy waste while considering his comfort. A case describes the energy management scenario of a building for one day. The actions retained in the system case base are the actions effectively carried out by the building occupant, so there is no guarantee that they are actions that generate satisfactory effects for the occupant. For this reason, the system is provided with a function to evaluate the performance of the actions stored in the case base, allowing to label the corresponding cases with the appropriate labels.

2.1 Founding notions and notations about CBR Approach

The memory of a CBR system is made of a set of source cases C_{sr} which constitute a case base CB.

2.1.1 CASE DESCRIPTION

Let \mathbb{C} , \mathbb{A} , and \mathbb{E} be three mutually disjoint sets. A case is a triplet $(\mathscr{C}, \mathcal{A}, \mathcal{E}) \in \mathbb{C} \times \mathbb{A} \times \mathbb{E}$ where:

- *C* is an element of the context domain ℂ, i.e., the imposed elements of the problem over which one cannot exert control. For instance, in a CBR-based disease treatment system, the context data can be the different physiological measures of the patient (blood pressure, glycemic rate, etc.).
- A is an element of the action domain A, i.e., elements that can be controlled to achieve the relevant outcomes. It represents the solution proposed by the system. For instance, the names and the protocol for administering the drugs prescribed in a CBR-based disease treatment system.
- *E* is an element of the effect domain 𝔼, i.e., elements describe the state of the system after applying action *A* to context *C*. For instance, the patient's physiological measures after the treatment.

A target context \mathscr{C}_{tg} is a context for which the CBR system tries to predict target actions \mathcal{A}_{tg} to generate target effects \mathcal{E}_{tg} and thus elaborate a target case C_{tg} . Formally, the resolution of a problem in the CBR paradigm is defined by Equation equation 1.

$$CBR \text{ system: } (CB, \mathscr{C}_{tg}) \longmapsto \mathcal{A}_{tg}$$

$$C_{tg} \triangleq (\mathscr{C}_{ta}, \mathcal{A}_{ta}, \mathcal{E}_{tg})$$
(1)

With CB – the case base.

2.1.2 RETRIEVING AND ADAPTATION

(

A full presentation of the reasoning process is beyond the focus of this paper, but due to the particular connection between adaptation and retrieving knowledge, it is usually necessary to present the adaptation process in conjunction with the retrieval process. Indeed, the reasoning process modeled by Equation equation 1 is made up of two steps.

• given a threshold σ for the distance between the context variables of the source cases and the target context, the retrieval process consists of identifying the source cases having a context similar to the target context. The profile of the retrieval function is given in Equation equation 2.

Retrieve:
$$\mathscr{C}_{tq} \longmapsto \{ \forall C_{sr} \in CB/Distance(\mathscr{C}_{tq}, \mathscr{C}_{sr}) \le \sigma \} = \mathcal{S}_{\mathcal{C}_{tq}}$$
 (2)

Where $Distance(\mathscr{C}_{tg}, \mathscr{C}_{sr})$ – a metric that computes the distance between the context variables \mathscr{C}_{tq} of the target case C_{tq} and the context variables \mathscr{C}_{sr} of source cases C_{sr} .

No constraints are imposed on the type of distance to use since it permits handling the context variables. For instance, the Minkowski metric can be used to calculate the context distance in a CBR-based EMS since the context variables are real values.

• since the source contexts usually do not match the target context, it is required to define a function to adapt the source actions to satisfy the requirements of the target context. The profile of the adaptation function is defined by the Formula equation 3.

Adaptation:
$$\forall C_{sr} \triangleq (\mathscr{C}_{sr}, \mathcal{A}_{sr}, \mathcal{E}_{sr}) \in \mathcal{S}_{\mathcal{C}_{tg}},$$

 $(\{(\mathscr{C}_{sr}, \mathcal{A}_{sr}, \mathcal{E}_{sr})\}, \mathscr{C}_{tg}) \longmapsto \mathcal{A}_{tg}$
(3)

Where $S_{C_{ta}}$ – the set of similar source cases as defined by equation equation 2.

Note that Equation equation 3 does not impose any constraints on the number of similar cases considered in the adaptation process, thus we are dealing with a compositional adaptation (whose single case adaptation is a particular case), where solutions from several source cases are combined to yield a target solution. Indeed, the experiment indicated that retaining a single case often gives less accurate results Sizov et al. (2016). This is explained by the fact that frequently only a part of the problem of the similar source case is relevant for the target problem, which makes the task of adaptation complicated (if not impossible).



Figure 1: Artificial potential field.

Figure 2: CBR attractive force

2.2 Collisionless path planning

Robot path planning study focus on the path planning of an autonomous robot moving in an unknown environment, i.e., guide the robot in its displacement from an initial position to a target position by calculating the optimal but moreover the safest path, i.e., avoiding obstacles that can occur along the path towards the target.

Several approaches were proposed to tackle this challenge, in particular, the artificial potential field approach originally proposed in Khatib (1985) is extensively adopted in robot guidance. The artificial potential field approach can cope with the reality of the current environment of the robot displacement by considering both the objectives to be reached and the obstacles to be avoided while moving. The key idea of this approach is to consider the robot as a point evolving in a 2-dimensional space (in the basic scenario) subject to the field influences of targets to reach and obstacles to avoid. Consequently, the robot is subjected to two kinds of forces, including an attractive one \mathbb{F}_{at} generated by targets and a repulsive one \mathbb{F}_{rp} generated by obstacles to move the robot further away.

Whereas repulsive forces are disproportional to the distance between the robot and the obstacles, i.e. they are strongest close to the obstacles and are less influential at distance, attractive forces are proportional to the distance between the target and the robot. The combined (total) of all the forces $\overrightarrow{\mathbb{F}} = \overrightarrow{\mathbb{F}_{at}} + \overrightarrow{\mathbb{F}_{rp}}$ applied to the robot defines the movement direction of the robot and its speed whilst avoiding collisions with obstacles. For the sake of simplification, the principle of this method for a robot traveling in a 2-dimensions environment is depicted in Figure 1.

3 REASONING FROM SUCCESSFUL AND FAILED CASES

3.1 PROBLEM FORMALIZATION

The adaptation problem considering failed and successful cases can be formalized as follows. Given the following observations:

• the case base CB is divided into two partitions of failed cases CB_{-} and successful cases CB_{+} . So,

$$CB = CB_{-} \cup CB_{+}$$

• by misuse of language, we refer to a target case as the elements of a target context for which we are looking for a solution. The case structure is not completely defined as the elements representing the actions and therefore, the effects are unknown.

Find a solution for a target case (thus under construction) is to infer, from source cases having similar context, a set of target actions that best satisfy the target context, which leads to the definition of the target effects, and thus to building an effective case containing the three elements: context, actions, and effects.

Similar source cases should be handled differently depending on whether they are failed (member of CB_{-}) or successful (member of CB_{+}) and on their degree of similarity to the target case. The method to be proposed should provide mechanisms to move towards the solutions of successful similar source cases and away from failed similar source cases while taking into account that the closer the source case to the target case the more influence its solution has on the target solution.

3.2 PRINCIPLE

The principle of our approach to considering failed cases in the adaptation process is inspired by navigation algorithms originating from the literature on the programming of autonomous robots, in particular, based on the artificial potential field presented in Section 2.2.

Before describing the details of our approach in the next section, to ensure the successful implementation of an artificial potential field-like concept in the context of this work, some assumptions are formulated:

• while the labeling process falls outside the scope of this study, we assume that previous experiences are already labeled as successful or failed cases. Furthermore, we suppose that the CBR system is given a quality function Q which scores the efficacy of the actions applied to the context. The highest scores are the best. This implicitly defines a threshold value $\mathcal{P}_s^{\mathcal{E}_i}$ for each effect feature \mathcal{E}_i according to Equation equation 4.

$$\forall C_i \in CB , \ \mathcal{Q} : \mathcal{E}_i \longmapsto \mathbb{R}$$
$$\mathcal{L}(C_i) = \begin{cases} C_i + & \text{if } \mathcal{Q}(\mathcal{E}_i) \geq \mathcal{P}_s^{\mathcal{E}_i} , \ \forall \mathcal{E}_i \in \mathbb{E} \\ C_i - & \text{otherwise.} \end{cases}$$
(4)

With \mathcal{L} – the labeling function, CB – the case base, \mathcal{E}_i – an effect feature of case C_i .

• classical CBR methods retrieve a defined number of neighboring cases from the case base CB regardless of an optimal number of similar ones regarding the target case. This KNN-like approach poses some issues since the target cases do not necessarily have the same number of similar neighbors, while some target cases should have more similar cases, others less. Furthermore, the configuration where much more source cases with equal distance from a target case than the predefined number, must be handled. In this work, we assume the existence of a retrieval approach that adjusts the number of source cases similar to the target case C_{tr} by dynamically defining a similarity threshold $\sigma_{C_{tr}}$ for the context distance between C_{tr} and the neighboring source cases. For instance, the work presented in Boulmaiz et al. (2021a), provides a method to define this threshold by combining a statistical approach and a genetic algorithm.

The key idea of the approach proposed in this work is to map the type of source cases available in the case base, i.e., successful and failed cases, to the type of objects handled in the context of robot moving, i.e., target and obstacles. Therefore, failed cases are assimilated into obstacles and successful cases into targets. While cases $C_i + \in S_{C_{tg}}$ with good performances should generate an attractive force \mathbb{F}_{at} that pulls the target solution towards them, the bad cases $C_i - \in S_{C_{tg}}$ should produce a repulsive force \mathbb{F}_{rp} that pushes away the solution from them.

The successful and failed source cases are considered to be sources for generating a potential field representing the properties of the target solution. As in the robotic potential field method, the CBR potential field is still composed of two fields. Regarding the attractive potential field, an attractive force is produced from the target solution to the source solutions of the successful cases by the configuration of the latter, which allows to pull the target solution towards the solutions of these cases.

To illustrate this concept, let's consider, for the sake of presentation, a system with domain knowledge containing only 2 action variables, the attractive potential field generated by any successful case looks like Figure 2, where at each point of the context space representing the target context, the force vectors are directed towards the successful source case. Concerning the repulsive potential field, a pushback force is generated by the configuration of the failed case towards the target solution, which allows to pull the target solution away from the solutions of these cases. Figure 3 depicts the CBR repulsive force in a similar configuration to the example illustrating the CBR attractive force.



Figure 3: CBR repulsive force

Figure 4: CBR total potential force

Ultimately, the configuration of the target solution, i.e., the position of the target solution in the space of solutions (actions), is determined by summing all repulsive and attractive forces generated by neighboring failed and successful cases respectively. For the simple case of only two neighbors, a successful case and a failed case, the total potential field has the shape shown in Figure 4.

3.3 LOCAL PREDICTION OF THE TARGET SOLUTION

Although we are inspired by the potential artificial field method, its application in the context of this work as applied in the robotics context does not permit determining the solution for many reasons:

- the potential total force in the robotic context depends exclusively on the distance between the goal/obstacles and the robot. In the CBR context, the magnitude of the attraction and repulsion forces are not dependent only on the distance between the target context and the neighboring source contexts but also on the performance of the neighboring source contexts.
- within the robotics context, unlike the attractive force, the magnitude of the repulsive force is at its highest value close to the obstacle and decreases proportionally when moving away from it. Within the context of CBR applications, the magnitude of the two forces should be proportional to the performance of the source solutions but disproportional to the distance between the source contexts and the target one.
- there is usually only one goal to reach in robotic applications, but in the case of a multi-goal environment, one looks for a path that goes through all these goals in sequential order by optimizing some criteria. for CBR systems, the aim is to combine the knowledge of all the neighboring source cases to infer the target solution.
- while the purpose of the robotic potential artificial field is to find the safe path to the goal, its purpose in the CBR application is to acquire new knowledge that guides the adaptation process in the construction of the target solution, i.e., to orient the reasoning process towards the most useful solutions (closest and best-performing cases) and away from the worst cases (farthest away or bad performance).

It is, therefore, necessary to adapt the approach of the artificial potential field to take into consideration the specificities of the CBR adaptation process. To do so, our approach defines the target solution (actions) \mathcal{A}_{tg} by the vectorial sum of all attractive forces $(\mathbb{F}_{at}^{C_i+}, \forall C_i+ \in \mathcal{S}_{\mathcal{C}_{tg}})$ and all repulsive forces $(\mathbb{F}_{rp}^{C_i-}, \forall C_i- \in \mathcal{S}_{\mathcal{C}_{tg}})$ as described in equation equation 5.

$$\sum_{C_i} \mathbb{F}^{C_i} \overrightarrow{\mathcal{A}_{tg} \mathcal{A}_i} = \sum_{C_i+} \mathbb{F}^{C_i+}_{at} \overrightarrow{\mathcal{A}_{tg} \mathcal{A}_{C_i+}} + \sum_{C_i-} \mathbb{F}^{C_i-}_{rp} \overrightarrow{\mathcal{A}_{tg} \mathcal{A}_{C_i-}} = 0$$
(5)

As already mentioned earlier, the magnitude of the repulsion and attraction forces depends both on the distance of the target context from the context of the similar source case and on the performance of the latter. From Equation equation 5, the metric \mathbb{F}^{C_i} defines the magnitude and direction of the associated force to the case C_i . We propose in Equation equation 6 a formula to estimate its value.

$$\forall C_i \in \mathcal{S}_{\mathcal{C}_{tg}},$$

$$\mathbb{F}^{C_i} = \begin{cases} \left(1 - \frac{\mathcal{D}_C(C_{tg}, C_i)}{\sigma_{C_{tg}}}\right) \times (\mathcal{Q}_i - \mathcal{P}_s) & \text{if } \mathcal{Q}_i \neq \mathcal{P}_s \\\\ 1 - \frac{\mathcal{D}_C(C_{tg}, C_i)}{\sigma_{C_{tg}}} & \text{else} \end{cases}$$
(6)

With $\sigma_{C_{tg}}$ – the context distance threshold, Q_i – the performance of the case C_i , \mathcal{P}_s – the performance threshold, $\mathcal{D}_C(C_{tg}, C_i)$ – the context distance between C_{tg} and its neighbor C_i .

From Equation equation 6, one can observe that whatever the type of force, its magnitude progressively decreases at the expense of an increasing context distance until it becomes null when the context distance equals the similarity threshold $\sigma_{C_{tg}}$. Besides defining the magnitude of the force, the operand $Q_i - P_s$ specifies the type of the force. When $Q_i \ge P_s$, then $\mathbb{F}^{C_i} \ge 0$, and the case C_i generates an attractive force else, it should be a repulsive force.

In this manner, the actions to be proposed A_{tq} have to satisfy:

$$\mathcal{A}_{tg} = \frac{1}{\sum_{C_i} \mathbb{F}^{C_i}} \sum_{C_i} \mathbb{F}^{C_i} \mathcal{A}_i , \ \forall C_i \in \mathcal{S}_{\mathcal{C}_{tg}}$$
(7)

Where $S_{\mathcal{C}_{tg}}$ – the set of neighboring cases to the target case C_{tg} .

4 EVALUATION

The objective of the evaluation is twofold, i) to study the potential impact of considering both failed and successful cases on improving the performance of the CBR system. ii) to assess the performance of the artificial potential field approach, this is referred to as *CBR-APF* in the following, compared to other adaptation approaches. To do so, several baselines are considered:

- 1. the approach proposed in Boulmaiz et al. (2022), denoted *CBR-S* in the following, exploits failed and successful cases but with a null adaptation process as the latter consists in making a vote among the similar cases solutions to select the solution with the best performance (maximizes the quality function) by applying it directly to the target case. pour tester la fait davoir plusieurs solution sources
- 2. a standard barycentric approach that combines solutions from the set of successful and failed similar source cases, noted *CBR-B* hereafter. tester le benefice des forces artificiaelle
- 3. a modified variant denoted *CBR-P* of our approach is tested, it considers only positive cases and thus uses only attractive forces. The objective is to illustrate the advantage of considering both negative and positive cases w.r.t only positive cases.
- 4. the approach proposed in Patterson et al. (2002) is used as a further baseline. This approach referred to as *CBR-R*, is based on a KNN approach to select similar source cases from which a generalized case is generated. Similar cases are used also to train a linear regression model, which is applied to the generalized case to predict the target case solution.

The performances of all approaches are compared to the ones of the actions recorded in the case base (actions effectively performed by the user without assistance, that are denoted *CBR-U* in the following) according to three measures: performance enhancement rate (PER), approach efficiency rate (APR), and effect quality rate (EQR). More details on the semantics and computation of these metrics are given in Appendix B.6.

4.1 EXPERIMENTAL SETUP

As mentioned in Section 2, the approach is implemented in an EMS whose objective is to make the user aware of the impact of his actions on the energy use in a building. Concretely, the EMS proposes to the occupant a series of actions to improve the comfort while consuming less energy. More precisely, given the weather forecast of a future day, the EMS looks to identify in its memory the past days with the same context to suggest the best action schedule to enhance the user's comfort for less energy consumption.

To evaluate our approach, we perform an initial experiment using real data collected from an university office (see Section B.2 for further information). The generating process of the case base is detailed in Section B.4.1. Unfortunately, after the pre-processing step described in Section B.3, the size of the case base retained is relatively small (98 cases). To perform a large-scale validation, we conduct a second experiment using semi-synthetic data generated from real data used in the initial experiment. More details on the semi-synthetic data generation process can be found in Section B.4.2.However, in the following, Due to space constraints, we have chosen to discuss only the results of the large-scale validation. Other results for the case base generated exclusively from real data can be found in Section D.

4.2 **Results and analyse**

Whatever the adaptation process approach applied in a CBR system, its performance depends partially on the retrieval process. Analyze the latter goes beyond the scope of the present paper, we detail in Appendix C the procedure applied to extracting similar source cases. It follows that each target case (test case) has at least one similar source case.

Table 1 summarizes the results of the 5-fold cross-validation of our approach against the four baselines considered. Some important findings from this experiment are:

- while the value of the EQR metric (see Appendix 4 for the exact definition) corresponds to the value of APR for the CBR-S, CBR-B, and CBR-APF approaches, the APR value is less than that of EQR for the CBR-P and CBR-R approaches, this is due to the ability of the first three approaches to computing a solution even with a similar set of cases consisting exclusively of failed cases.
- our CBR-APF approach is clearly better in performance than all other baselines with also better RPA and EQR, regardless of the test set.
- the number of similar source cases has a significant influence on the quality of the adaptation process, a compositional adaptation (which uses several similar source cases) systematically gives a better PER, as illustrated by the comparison between PERs of CBR-APF which is a compositional approach and CBR-S which uses a single similar case.
- attraction and repulsion forces have an important impact on the results of the adaptation process. Given the same number of similar cases, by using these forces, our CBR-APF approach outperforms the CBR-B baseline, which does not use them. CBR-APF is 1.64% times more performing than CBR-B regarding the improvement of the cases performances (global PER = 33.49% versus 20.36%) and 1.61% times more efficient according to the number of cases for which it manages to find a solution (CBR-APF improves the performance of the solutions proposed by the user without assistance for 99.98% of cases against 61.87% for CBR-B).
- using failed cases in case-based reasoning significantly influences the performance of a CBR system. By exploiting both successful and failed cases, the system improves the results of the reasoning process. Comparing the performance of the CBR-APF approach with that of the CBR-P and CBR-R approaches (both do not use failed cases in their reasoning), the EQR results (ratio of the number of cases whose performance is improved to the number of cases whose solutions are found with either improved or degraded performance) show that the CBR-APF approach outperforms the other baselines. Our CBR-APF approach is more than three times more efficient than CBR-R and more than 1.5 times more than CBR-P in improving the performance of test cases.

The validation results on a smaller real dataset are quite similar (see Appendix 5). This strengthens the initial hypothesis motivating this work, which is using failed cases jointly with successful cases enhances the adaptation process.

_	TEST SET		S1			S2	-		S3			S4			S5		(GLOBA	L
APPROACH			METRICS		N N	IETRIC	s	N	IETRIC	s	N	IETRIC	s	1	METRIC	s	l N	IETRIC	s
		PER (%)	APR(%)	EQR(%)	PER	APR	EQR	PER	APR	EQR	PER	APR	EQR	PER	APR	EQR	PER	APR	EQR
CBR - S		16.73	59.13	59.13	17.85	48.57	48.57	19.53	60.12	60.12	20.48	56.07	56.07	18.79	64.48	64.48	18.68	57.67	57.67
CBR - B		18.27	57.51	57.51	15.36	63.90	63.90	22.85	59.69	59.69	24.23	65.52	65.52	21.10	662.71	62.71	20.36	61.87	61.87
CBR - P		22.62	42.26	57.10	18.54	48.85	63.71	20.14	50.21	60.10	22.48	52.92	70.19	23.47	39.86	60.09	21.45	46;82	62,24
CBR - R		-2.56	32.18	49.75	9.12	29.89	51.19	14.71	43.07	64.24	17.45	39.52	57.74	12.04	41.26	62.84	10.15	37.18	57.15
CBR - APF		34.68	100	100	28.85	99.76	99.76	33.91	100	100	31.27	100	100	38.73	99.88	99.88	33.49	99.92	99.92

 Table 1: Summery of results on synthetic dataset

5 CONCLUSION

This paper proposed a new approach to the adaptation process in the CBR paradigm by looking at both failed and successful source cases instead of the traditional practice of considering only successful source case. We found inspiration in the studies on planning safe paths for a robot moving in an unknown environment. The concept is that both successful and failed cases generate attraction and repulsion forces respectively on a likely barycentric solution to drive the reasoning towards the best performing solutions and away from the failed ones. The experimentation of this approach in the context of an energy management system showed a significant improvement in the system performance by considering both successful cases and failed ones.

Compared to the only existing work Lieber & Nauer (2021) based on the same logic as ours (see Section 1), the advantages of our approach are as follows:

- the approach we have proposed does not impose a Boolean notation to represent the components of a case, which is the case in Lieber & Nauer (2021) that requires the transformation of the problem and solution attributes into a Boolean representation.
- our approach proposes a fully domain-independent adaptation approach that does not require any expert/user intervention, while the approach presented in Lieber & Nauer (2021) is a domain-independent task where an expert/user is solicited to establish adaptation rules between each pair of source cases. Indeed, the adaptation process adopted in Lieber & Nauer (2021) is based on the well-known difference between cases approach that consists in establishing adaptation rules by studying the influence of the differences between the problems of two source cases on the difference between their solutions.

Unfortunately, an obvious negative criticism that can be formulated against our approach is it does not guarantee to find a successful solution when all similar source cases are failed cases. Indeed, in this particular case, our approach should propose a solution (actions plan) that is certainly far from the failed solutions of similar source cases but this is not sufficient to grant that the proposed solution is successful. However, despite if it is a failed solution, this is interesting from a knowledge acquisition perspective since this failed case might be used in the future to compute solutions for other target cases.

In this work we have developed and evaluated an approach considering the whole set of successful and failed cases, it would be interesting to perform a deeper evaluation taking into account the number of neighboring successful and failed cases considering only the n cases with the best performances and the m cases with the worst performances. Another line of future research for this work would be to explore the possible impact of a failed case on the domain ontology (if any). It could be useful to suggest new necessary conditions to add to the domain ontology that would avoid the reappearance of such a negative case in the future.

REFERENCES

- The International Energy Agency. Building energyuse in chinatransforming constructionand influencingconsumption to 2050. Technical report, Tsinghua University Building Energy Research Center, 2015.
- ASHRAE (ed.). ASHRAE Standard Thermal Environmental Conditions for Human Occupancy. American Society of Heating, Refrigerating and Air-Conditioning Engineers., Atlanta,USA, 1992.

- ASHRAE (ed.). Indoor air quality guide: best practices for design, construction, and commissioning. American Society of Heating, Refrigerating and Air-Conditioning Engineers., Atlanta, USA, 2009.
- Fadi Badra, Julien Cojan, Amélie Cordier, Jean Lieber, Thomas Meilender, Alain Mille, Pascal Molli, Emmanuel Nauer, Amedeo Napoli, Hala Skaf-Molli, and Yannick Toussaint. Knowledge acquisition and discovery for the textual case-based cooking system wikitaaable. 8th International Conference on Case-Based Reasoning-ICCBR 2009, 2009.
- Fateh Boulmaiz, Amr Alzouhri Alyafi, Stephane Ploix, and Patrick Reignier. Optimizing occupant actions to enhance his comfort while reducing energy demand in buildings. In *11th IEEE IDAACS*, 2021a.
- Fateh Boulmaiz, Stephane Ploix, and Patrick Reignier. A data-driven approach for guiding the occupant's actions to achieve better comfort in buildings. In 2021 IEEE 33rd ICTAI, pp. 463– 468, 2021b. doi: 10.1109/ICTAI52525.2021.00075.
- Fateh Boulmaiz, Patrick Reignier, and Stephane Ploix. An occupant-centered approach to improve both his comfort and the energy efficiency of the building. *Knowledge-Based Systems*, 249:108970, 2022. ISSN 0950-7051. doi: https://doi.org/10.1016/j.knosys. 2022.108970. URL https://www.sciencedirect.com/science/article/pii/ S0950705122004701.
- Amélie Cordier, Béatrice Fuchs, Léonardo Lana de Carvalho, Jean Lieber, and Alain Mille. Opportunistic acquisition of adaptation knowledge and cases the IakA approach. In *Lecture Notes in Computer Science*, pp. 150–164. Springer Berlin Heidelberg. doi: 10.1007/978-3-540-85502-6_10.
- Amélie Cordier, Béatrice Fuchs, Léonardo Lana de Carvalho, Jean Lieber, and Alain Mille. Opportunistic acquisition of adaptation knowledge and cases - the iaka approach. In *ECCBR*, 2008.
- Susan Craw, Jacek Jarmulak, and Ray Rowe. Learning and applying case-based adaptation knowledge. In David W. Aha and Ian Watson (eds.), *Case-Based Reasoning Research and Development*, pp. 131–145, Berlin, Heidelberg, 2001. Springer Berlin Heidelberg.
- Susan Craw, Nirmalie Wiratunga, and Ray C. Rowe. Learning adaptation knowledge to improve case-based reasoning. *Artificial Intelligence*, 170(16):1175–1192, 2006. ISSN 0004-3702. doi: https://doi.org/10.1016/j.artint.2006.09.001. URL https://www.sciencedirect.com/science/article/pii/S0004370206000798.
- Elham Delzendeh, Song Wu, Angela Lee, and Ying Zhou. The impact of occupants' behaviours on building energy analysis: A research review. *Renewable and Sustainable Energy Reviews*, 80:1061–1071, 2017. ISSN 1364-0321. doi: https://doi.org/10.1016/j.rser. 2017.05.264. URL https://www.sciencedirect.com/science/article/pii/ S1364032117309061.
- Belén Díaz-Agudo and Pedro A. González-Calero. An architecture for knowledge intensive CBR systems. In *Lecture Notes in Computer Science*, pp. 37–48. Springer Berlin Heidelberg, 2000. doi: 10.1007/3-540-44527-7_5.
- Simona D'Oca and Tianzhen Hong. A data-mining approach to discover patterns of window opening and closing behavior in offices. *Building and Environment*, 82:726–739, 2014. ISSN 0360-1323. doi: https://doi.org/10.1016/j.buildenv.2014.10.021. URL https://www.sciencedirect.com/science/article/pii/S0360132314003424.
- International Organization for Standardization. Building environment design indoor air quality methods of expressing the quality of indoor air for human occupancy, 2008.
- Béatrice Fuchs, Jean Lieber, Alain Mille, and Amedeo Napoli. An algorithm for adaptation in case-based reasoning. pp. 45–49, 01 2000.
- Nadezhda Govedarova, Stanimir Stojanov, and Ivan P. Popchev. An ontology based cbr architecture for knowledge management in bulchino catalogue. In *CompSysTech*, 2008.

- CSA Group. Z412-17 Office ergonomics An application standard for workplace ergonomics. 2017.
- Kathleen Hanney and Mark T. Keane. Learning adaptation rules from a case-base. In Ian Smith and Boi Faltings (eds.), Advances in Case-Based Reasoning, pp. 179–192, Berlin, Heidelberg, 1996. Springer Berlin Heidelberg. ISBN 978-3-540-49568-0.
- Kathleen Hanney and Mark T. Keane. The adaptation knowledge bottleneck: How to ease it by learning from cases. In *Case-Based Reasoning Research and Development*, pp. 359–370. Springer Berlin Heidelberg, 1997. doi: 10.1007/3-540-63233-6_506.
- Vahid Jalali and David Leake. Scaling up ensemble of adaptations for classification by approximate nearest neighbor retrieval. In David W. Aha and Jean Lieber (eds.), *Case-Based Reasoning Research and Development*, pp. 154–169, Cham, 2017. Springer International Publishing.
- O. Khatib. Real-time obstacle avoidance for manipulators and mobile robots. In *Proceedings*. *IEEE International Conference on Robotics and Automation*, 1985. doi: 10.1109/ROBOT.1985. 1087247.
- Sergei Kuznetsov. Machine learning on the basis of formal concept analysis. *Automation and Remote Control*, 62:1543–1564, 10 2001. doi: 10.1023/A:1012435612567.
- Jared Langevin, Jin Wen, and Patrick L. Gurian. Including occupants in building performance simulation: Integration of an agent-based occupant behavior algorithm with energyplus. 2014 ASHRAE/IBPSA-USA Building Simulation Conference, September 2014.
- Jean Lieber and Emmanuel Nauer. Adaptation knowledge discovery using positive and negative cases. In Antonio A. Sánchez-Ruiz and Michael W. Floyd (eds.), *Case-Based Reasoning Research and Development*, pp. 140–155, Cham, 2021. Springer International Publishing. ISBN 978-3-030-86957-1.
- Neil McDonnell and Pádraig Cunningham. A knowledge-light approach to regression using casebased reasoning. In Thomas R. Roth-Berghofer, Mehmet H. Göker, and H. Altay Güvenir (eds.), *Advances in Case-Based Reasoning*, pp. 91–105, Berlin, Heidelberg, 2006. Springer Berlin Heidelberg. ISBN 978-3-540-36846-5.
- David McSherry. An adaptation heuristic for case-based estimation. In Barry Smyth and Pádraig Cunningham (eds.), Advances in Case-Based Reasoning, pp. 184–195, Berlin, Heidelberg, 1998. Springer Berlin Heidelberg.
- Mirjam Minor and Lutz Marx. Case-based Reasoning for Inert Systems in Building Energy Management. In Proc. ICCBR 2017, pp. 200–211. Springer, 2017.
- Tom M. Mitchell. Generalization as search. *Artificial Intelligence*, 18(2):203–226, 1982. ISSN 0004-3702. doi: https://doi.org/10.1016/0004-3702(82)90040-6.
- Rudradeb Mitra and Jayanta Basak. Methods of case adaptation: A survey. *International Journal of Intelligent Systems*, 20, 2005.
- Ligia Moga and I. Moga. Building design influence on the energy performance. *Journal of Applied Engineering Sciences*, 5(1):37–46, 2015. doi: doi:10.1515/jaes-2015-0005. URL https://doi.org/10.1515/jaes-2015-0005.
- Monalisa Pal, Amr Alzouhri Alyafi, Sanghamitra Bandyopadhyay, Stéphane Ploix, and Patrick Reignier. Enhancing comfort of occupants in energy buildings. In Samarjit Kar, Ujjwal Maulik, and Xiang Li (eds.), *Operations Research and Optimization*, pp. 133–144, Singapore, 2018. Springer Singapore. ISBN 978-981-10-7814-9.
- David W. Patterson, Niall Rooney, and Mykola Galushka. A regression based adaptation strategy for case-based reasoning. In *AAAI/IAAI*, 2002.
- Sanja Petrovic, Gulmira Khussainova, and Rupa Jagannathan. Knowledge-light adaptation approaches in case-based reasoning for radiotherapy treatment planning. *Artificial intelligence in medicine*, 68:17–28, 2016.

- Anna Laura Pisello and Francesco Asdrubali. Human-based energy retrofits in residential buildings: A cost-effective alternative to traditional physical strategies. *Applied Energy*, 133:224–235, 2014. ISSN 0306-2619. doi: https://doi.org/10.1016/j.apenergy.2014.07.049. URL https://www. sciencedirect.com/science/article/pii/S0306261914007314.
- Claudio Policastro, Andre Carvalho, and Alexandre Delbem. A hybrid case adaptation approach for case-based reasoning. *Applied Intelligence*, 28:101–119, 04 2008. doi: 10.1007/s10489-007-0044-4.
- Paula Rocha, Michal Kaut, and Afzal S. Siddiqui. Energy-efficient building retrofits: An assessment of regulatory proposals under uncertainty. *Energy*, 101:278–287, apr 2016. doi: 10.1016/j.energy. 2016.01.037.
- Karin Schakib-Ekbatan, Fatma Zehra Çakıcı, Marcel Schweiker, and Andreas Wagner. Does the occupant behavior match the energy concept of the building? analysis of a german naturally ventilated office building. *Building and Environment*, 84:142–150, 2015. ISSN 0360-1323. doi: https://doi.org/10.1016/j.buildenv.2014.10.018. URL https://www.sciencedirect.com/science/article/pii/S0360132314003394.
- Gleb Sizov, Pinar Öztürk, and Erwin Marsi. Compositional adaptation of explanations in textual case-based reasoning. In *Case-Based Reasoning Research and Development. ICCBR 2016*, pp. 387–401. Springer International Publishing, 2016. doi: 10.1007/978-3-319-47096-2.26.
- Wolfgang Wilke, Ivo Vollrath, Klaus-Dieter Althoff, and Ralph Bergmann. A framework for learning adaptation knowledge based on knowledge light approaches. In *Fifth German Workshop on Case-BasedReasoning*, pp. 235–242, 03 1997.
- Tzu-Tsung Wong. Performance evaluation of classification algorithms by k-fold and leave-oneout cross validation. *Pattern Recognition*, 48(9):2839–2846, 2015. ISSN 0031-3203. doi: https://doi.org/10.1016/j.patcog.2015.03.009. URL https://www.sciencedirect.com/ science/article/pii/S0031320315000989.
- Geun Young Yun and Koen Steemers. Behavioural, physical and socio-economic factors in household cooling energy consumption. *Applied Energy*, 88(6):2191-2200, 2011. ISSN 0306-2619. doi: https://doi.org/10.1016/j.apenergy.2011.01.010. URL https://www.sciencedirect.com/science/article/pii/S0306261911000134.
- Bingqing Zhang, Xiaodong Li, and Yimin Zhu. A cbr-based decision-making model for supporting the intelligent energy-efficient design of the exterior envelope of public and commercial buildings. *Energy and Buildings*, 231:110625, 2021. ISSN 0378-7788. doi: https://doi.org/10.1016/j.enbuild.2020.110625.
- Yue Zhao, Zain Nasrullah, and Zheng Li. Pyod: A python toolbox for scalable outlier detection. *Journal of Machine Learning Research*, 20(96):1–7, 2019. URL http://jmlr.org/ papers/v20/19-011.html.

A RELATED WORK

A.1 INDEPENDENT DOMAIN ADAPTATION APPROACHES

The authors in Mitra & Basak (2005) reviewed existing adaptation approaches and classified them according to three criteria: the need for domain knowledge in the specification of adaptation rules, generating techniques for adaptation methods, and adaptation knowledge type.

A classical design challenge for the case-based reasoning paradigm is to develop approaches independently as much as possible from the application domain to acquire the relevant knowledge for the adaptation process. The adaptation process is so difficult that most CBR systems are developed to leave it to domain experts. Unfortunately, manual adaption rule generation needs deep domain knowledge, making this process a difficult, costly, and inaccurate task in many domains. To tackle this issue, "lightweight" approaches have been elaborated to generate adaptation knowledge with minimal or no expert intervention. A widespread approach to generating adaptation rules is the application of machine learning techniques. The authors in Patterson et al. (2002) describe how to build adaptation rules by applying a linear regression model to the differences between the similar source cases to the target case. The method presented in Craw et al. (2006) proposes to use decision trees for acquiring adaptation knowledge. An approach consisting in splitting the case base into a small set of reference cases used to extract cases to produce adaptation rules from is suggested in Craw et al. (2001). In the study Policastro et al. (2008), a two-component algorithm is developed to automatically perform the adaptation process. Estimators based on Multi Layer Perception (MLP) neural network, M5 regression tree learner, and Support Vector Machine (SVM) technique are applied to produce a data set used by the combiner to generate adaptation rules.

Alternative approaches based on data mining techniques have been introduced to minimize the intervention of an expert in the adaptation rule elaboration process through an interactive process between the expert and the CBR system, as presented in Badra et al. (2009); Cordier et al..

The foundation of a particularly rich stream of research on domain-independent adaptation issue is based on the case difference heuristic approach first introduced in Hanney & Keane (1997). The rationale of this approach is to compare two source cases by attributing the differences between the solutions of the two cases to the differences in their problems. This generates an adaptation rule applicable to a source case having the same differences with the problem of a target case by adjusting the solution of the solution of the solution of the target case. The authors in Jalali & Leake (2017) combine the case difference heuristic approach and a linear regression model to acquire adaptation knowledge for systems described by non-numerical features. The variant reported in McSherry (1998) generates an adaptation rule using three source cases instead of a single source case. The principle of the differential adaptation presented Fuchs et al. (2000) consists of evaluating the variations between the descriptors (features) of the similar source case and the target case and their impact on the variation of the solution descriptors by defining a set of similarity and dissimilarity rules.

A.2 REASONING FROM FAILED AND SUCCESSFUL CASES

To our knowledge, the only study investigating so far this concept was reported in Lieber & Nauer (2021), where the authors proposed a three-phase approach. Firstly, an approach is described to transform a description of a case originally encoded by non-Boolean features into a Boolean encoding. Secondly, the generation of adaptation rules by applying the case difference heuristic approach Hanney & Keane (1996) between positive cases. The choice of the adaptation rule to apply is performed using symbolic data mining algorithms, namely the frequent closed set approach (FCI). Lastly, the use of only positive cases could generate too general adaptation rules whose application is liable to yield erroneous solutions, to overcome this issue, the negative (failed) cases are used to filter the too general adaptation rules by applying version space Mitchell (1982) and formal concept analysis Kuznetsov (2001) techniques.

Although the evaluation of this approach gives promising results, it suffers from some weaknesses:

- the formalism of the Boolean representation imposes to handle cases where the problem and its solution are represented by Boolean attributes or by attributes that are readily transformable into Boolean ones (for instance nominal data).
- learning adaptation knowledge is a domain-dependent task since it consists in comparing cases pairwise. Such an expert/user-based comparison has some drawbacks: (1) comparison may not be possible in complex applications characterized by many dependent features, such as EMS, (2) it is a time-consuming process requiring a high level of expertise, (3) it requires the availability of an expert every time the case base is updated with new case knowledge to update the adaptation rules.

B CASE STUDY: ENERGY MANAGEMENT IN BUILDING

This section provides details on the process of generating the case base used for evaluating our proposals.

B.1 MOTIVATION

Energy is a major economic and strategic issue for modern societies. Indeed, over the last few decades, the decline in natural energy resources (oil and gas) is coupled with a significant rise in consumption of energy and associated CO_2 emissions as a result of population growth and increasing comfort requirements of people. Given that buildings consume nearly 40% of primary energy in the European Union Rocha et al. (2016) and China Agency (2015), representing for the latter a growth of 37% between 2000 and 2012, and if this trend persists, it could grow by a further 70% by 2050, the reduction of energy consumption in buildings becomes a necessity.

Roughly speaking, building energy consumption can be regulated by two strategies:

- building material and appliance technology development strategy: This strategy includes approaches focusing on physical enhancements aimed at optimizing the performance of building envelopes Moga & Moga (2015) and appliances consuming energy Yun & Steemers (2011), such as heating, ventilation, and air conditioning systems.
- occupant behavior analysis-based strategy Pisello & Asdrubali (2014): Based on the finding that occupants form an integral part of the building energy behavior cycle Schakib-Ekbatan et al. (2015), it is necessary to analyze the dynamics of occupants' interactions with building systems to forecast energy use. The objective of this strategy is to give occupants an active role in reducing energy use and improving their comfort by providing them with energy information that motivates them to modify their decisions and behaviors in a greener direction.

Unlike the first strategy which is costly to perform physical retrofits, the second strategy is a promising alternative since it can be implemented, even in buildings already built, at zero cost by intelligently guiding the occupant's actions. Furthermore, the feeling of taking control of the factors influencing comfort yields a significant improvement in energy savings while avoiding the recognized discontent of occupants dismissed from the control. This highlights that a successful energy management strategy is closely tied to the interaction between the occupant and the building. Early studies D'Oca & Hong (2014); Delzendeh et al. (2017), conclude that basic actions, including windows and doors opening/closing, would reduce energy consumption. Nevertheless, the expected comfort may not be reached by such basic actions under extreme climatic conditions. To address this issue, the authors in Langevin et al. (2014) observed that supplementary heating/ cooling systems could improve the comfort of occupants and, therefore, aid in the energy management of buildings. Yet, few studies have been conducted to investigate how to include all of these occupant actions to achieve energy saving in buildings while maximizing occupant comfort.

Motivated by the above discussion, the aim is to provide a user-centered system by proposing an action plan to achieve maximum thermal and air quality satisfaction while expending minimal energy.

B.2 ENERGY DATASET

The evaluation is performed using real data collected from a workspace in a French university building. Numerous sensors are deployed in the workspace to collect data on the state of the indoor environment. We use a metheological service provider to collect data on the outdoor environment and the weather forecast. The set of data is represented by variables modeling the measured phenomena. Table 2 classifies the variables into one of the three components of a case according to the case structure adopted in Section 2.1.1.

Sensors generate data in a stream form at different time intervals, this experiment adopts an hourly granularity, and therefore the data streams are sampled each hour. The value considered for each feature at a given hour is the average data flow received during that hour, except for the two features window opening and door opening, where the values represent the fraction of the hour during which the window/door is fully opened. For instance, Ed=0.75 for the k^{th} hour corresponds to completely opening the door for 45 minutes (75% of the k hour).



B.3 PREPROCESSING

Like any data-driven approach, the data preprocessing stage is crucial to enhance the performance of the learning algorithms for both the retrieval and adaptation stages. Three types of preprocessing are applied to the database: date cleaning, data normalization and data filtering.

B.3.1 DATA CLEANING

Data acquisition from sensors installed in the office could induce errors in the collected data, for instance, missing values or outliers. So it is necessary to deal with these errors to guarantee a correct reasoning process. We use the AvgKNN class of the Python Outlier Detection framework (PyOD) Zhao et al. (2019) to detect outliers which are replaced by the average of the previous three nearest neighbors (based on the acquisition date) and the next three nearest neighbors.

B.3.2 FEATURES NORMALIZATION

Since the features presented in Table 2 describe different physical phenomena that could have a large range of values, features' values need to be normalized to eliminate the features' magnitudes to balance their contributions and compare their influences. We use MinMax normalization to scale features' values between 0 and 1.

B.3.3 HOURS FILTERING

As already mentioned, for this experiment, we study an office at the University of Grenoble, in France, where two scholars work between 8:00 am and 8:00 pm on working days. Consequently, we filter out hours when no one is at the office, i.e., weekends, and public holidays. This is motivated by the absence of actions during these time slots.

However, during weekdays, even if there is no one in the office during the night hours between 9:00 pm and 7:00 am (there are no actions), we consider that night hours in the reasoning process to account for the inertia of the building, i.e., the capacity of the physical envelope of the building to dynamically store energy. This phenomenon induces a time shift and smoothing effect on the various changes in the building's interior environment, avoiding sudden variations in temperature or air quality.

B.4 CASE BASE GENERATION AND TESTBED

Collected data from the different sensors and meteorological information are stored in a database starting in April 2015 and ending in October 2016. So, data are sampled hourly, a day is described by a 24-value vector for each variable and represents one case in the case base. After the preprocessing phase, a case base *CB* contains 98 cases (days) is considered for the experiment.

To validate our approach the case base *CB* is split in two disjoint sets: $CB = CB^{L} \cup CB^{T}$, the learning set CB^{L} which represents the source cases, and the test set CB^{T} which represents the target cases for which a solution is sought. However, to compensate for the problem of the limited number of real data and to give more reliability to the validation process, we have decided to use two datasets with different approaches for sampling the test and training data. The two datasets and the associated splitting process are described as follows:

B.4.1 REAL DATASET

In real dataset experiment, we exclusively use the real data corresponding to the case study described previously.

A common strategy used in the literature for splitting the data is the single holdout method, by randomly selecting a subset of the validation set (the case base) to be used as a test set, which often has about 10% to 30% of the available cases, the remaining cases forming the training set, which therefore includes about 90% to 70% of the cases. The holdout method is recommended when the validation data size is large enough to assure a valid test, unfortunately, in our experimentation the number of cases is limited, so we use the "leave-one-out cross-validation" (LOOCV) method, which is adopted for algorithms evaluation when only small number of instances are available Wong (2015).

The principle of the LOOCV validation is that one single case is considered as the test set ($|CB^{T}| = 1$), and all the other cases of the case base *CB* are considered as the learning set ($|CB^{L}| = 97$). This process is reiterated for each of the cases in the case base *CB*. In this way, the number of test cases is much larger than the traditional holdout approach. The results of this experiment are summarized in Appendix D.

B.4.2 SEMI-SYNTHETIC DATASET

To avoid deficiencies in available dataset, we are developing a tool to generate semi-synthetic data based on the real data already collected in the office. The semi-structured data generation consists in extending the real data set collected in the office. The principle of the generation approach is to synthesize the subset of action variables and to use real values for the context variables to generate effect variables. Precisely, the generating process includes three steps which are summarized as follows:

- context variables are generated from meteorological data provided by an online weather data provider provider for the period 01/01/1979 31/08/2022 for a total size of 15,948 days (|CB| = 15,948). Positive cases represent only 34.45% of the case base, i.e., the occupant failed to choose the optimal actions to achieve his comfort in 65% of the cases (see Table 3).
- to predict door and window opening/closing actions from the context variables of the 15,948 cases, we developed a dynamic Bayesian network-based tool that models occupants' actions by probabilistic cause-effect relations derived from expert knowledge as well as by conditional probabilities inferred from observations (real data collected in the office).
- to simulate the indoor conditions of the office (the effect variables) generated following the application of the actions to the corresponding contexts, a physical knowledge model of the office is elaborated. It consists of a set of mathematical equations that describe the energetic behavior and the air quality in the office.

For the validation process, we perform a 5-fold cross-validation where the original case base is initially randomly split into five equal-sized subsets: S1, 21, S3, S4, and S5. However, as the case base is moderately imbalanced with a ratio of 0.37 of positive cases to negative cases (see statistics in Table 3), the process of splitting the case base could induce a skewed distribution. To deal with this issue, we proceed as follows:

- the case base CB is split into two sets according to the performance of the cases, a set of positive cases CB₊ and a set of negative cases CB₋, i.e., CB = CB₊ ∪ CB₋.
- each of the sets CB_+ and CB_- is randomly split in 5 subsets of equal size: $S1_+$, $S2_+$, $S3_+$, $S4_+$, $S5_+$ and $S1_-$, $S2_-$, $S3_-$, $S4_-$, $S5_-$ respectively.
- each of the five final sets S1, S2 S3, S4, and S5 involved in the 5-fold cross-validation is formed by merging two sets taken randomly from sets {S1₊, S2₊, S3₊, S4₊, S5₊} and {S1₋, S2₋, S3₋, S4₋, S5₋} respectively. Table 3 provides more details on the distribution of positive and negative cases among the five sets used in the cross-validation.

Next, a single set is selected as a test set CB^T (target cases) while the remaining four sets are used as a learning set CB^L (source cases). The cross-validation procedure is performed five times, each

Tuble 5. Synthetic dutaset Statistics									
DATASET	C+ CASES	C- CASES	RATIO	TOTAL CASES	THRESHOLD σ	RECALL	PRECISION	F1-MEASURE	
S1	861	2,329	0.37%	3,190	1.13	77.82	69.94	73.67	
S2	861	2,329	0.37%	3,190	1.16	67.23	63.75	65.44	
S3	861	2,329	0.37%	3,190	1.13	72.80	70.17	71.46	
S4	861	2,329	0.37%	3,190	1.17	79.06	68.10	73.17	
S5	861	2,327	0.37%	3,188	1.15	63.96	76.90	69.83	
CB	4,305	11,643	0.37%	15,948	-	-	-	-	

 Table 3: Synthetic dataset Statistics

of the five sets being used once as a test set. The results of the metrics adopted to evaluate the performance are averaged to provide a final estimate.

B.5 QUALITY FUNCTION

In the EMS context, the performance function is related to the estimation of the user's satisfaction with the effects obtained as a result of the application of the actions proposed by the system. In this experiment, we consider two effect variables (see Table 2), indoor temperature and indoor CO_2 concentration.

The following values are established by experts in widely adopted standards (ASHRAE (1992) for the temperature comfort and for Standardization (2008) for the CO₂ concentration) in the residential and tertiary residence sector. The indoor temperature is considered good between the values 21°c and 23°c, it is acceptable within the intervals]23, 26] and [18, 21[, and it is bad outside these intervals. These values allow us to model the thermal dissatisfaction by the performance function presented in equation 8. The latter determine the temperature performance threshold $\mathcal{P}_s(T)$ that is fixed to 1, i.e., the temperature performance is considered bad when $\mathcal{Q}_T(t) > 1$.

Concerning the CO₂ concentration, the air quality is evaluated as good when the concentration of CO₂ is below 500 ppm, as acceptable within the interval [500, 1500], and otherwise it is bad. From these values, we can establish the Function equation 9 which models the performance function for CO₂ concentration. From these values, one defines the CO₂ performance threshold $\mathcal{P}_s(C)$ of 1 beyond which the performance in terms of CO₂ concentration is considered bad.

$$\mathcal{Q}_{T}^{h}(t) = \begin{cases} 0 & \text{if } t \in [21, 23] \\ \frac{t-23}{26-23} & \text{if } t > 23 \\ \frac{21-t}{21-18} & \text{if } t < 21 \end{cases}$$
(8)

$$Q_C^h(c) = \begin{cases} 0 & \text{if } c \le 500\\ \frac{c-500}{1500-500} & \text{if } c > 500 \end{cases}$$
(9)

$$Q^h = \frac{Q_T^h(t) + Q_C^h(c)}{2} \tag{10}$$

Global system performance at the *h*-hour Q^h is the average of the two performances Q_T^h and Q_C^h as described in Equation equation 10. A case is considered successful if the user is satisfied with both thermal comfort and air quality, otherwise, it is a failed case. This is formalized in Equation equation 11. Naturally, these thresholds are scalable and can be refined later to meet the user's comfort requirements.

Columns 2 and 3 of Table 3, give statistics concerning the number of successful cases and failed cases respectively for each test set. It can be noted from Table 4, which represents some statistics on the performance of the source cases in the training set, that the source solutions are performing well regarding air quality since the performance of these solutions is below the CO₂ performance threshold ($\mathcal{P}_s(C) = 1$).

$$\forall C_i \in CB,$$

$$\mathcal{L}(C_i) = \begin{cases} C_i + & \text{if } \mathcal{Q}_T^h(t) \leq \mathcal{P}_s(T) \land \mathcal{Q}_C^h(c) \leq \mathcal{P}_s(C), \forall h \in [0, 23] \\ C_i - & \text{otherwise.} \end{cases}$$
(11)

Table 4: Dissatisfaction Statistics							
DISSATISFACTION	MAX. VALUE	MIN. VALUE					
$\mathcal{Q}_s(T)$	2.3306	0.0098					
$\mathcal{Q}_s(C)$	0.0624	0.0					

With \mathcal{L} – the labeling function, CB - the case base, $\mathcal{P}_s(T)$ – the temperature threshold, $\mathcal{P}_s(C)$ – the CO₂ threshold.

A physical model Pal et al. (2018) of the workspace is used to simulate the effects of the actions proposed to the user to compare them to the actions effectively performed by the user. However, for the sake of consistency, the effects of the latter are also simulated by the physical model.

B.6 PERFORMANCE METRICS

We use the following metrics:

• performance enhancement rate (PER): The efficiency of the different baselines approaches (see Section 4) is evaluated by comparing, for each test case C_i , the average of the thermal performances Q_T^* , the air quality performances Q_C^* , and the global performance Q^* of the proposed actions to the corresponding values Q_T^r , Q_C^r , and Q^r of the actions already recorded in the case base. The performance enhancement \mathcal{H}_{C_i} related to the test case C_i , if any, is given by the equation 12. The global performance improvement of a baseline approach is calculated by averaging the performance improvements of all test cases.

$$PER_{C_i} = \frac{\mathcal{Q}^* - \mathcal{Q}^r}{\mathcal{Q}^r} \tag{12}$$

• *approach efficiency rate (APR)*: Global efficiency of a baseline approach is defined as the average of the ratio of the number of test cases whose performance is improved by applying the actions recommended by this approach to the total number of test cases.

$$APR = \frac{Z^+}{Z} \tag{13}$$

With $Z = |CB^T|$ – the set of test cases, $Z^+ = \{C_i \in CB_T \mid PER_{C_i} > 0\}$

• *effects quality rate (EQR)*: This measure is the average of the ratio between the number of test cases whose performance is improved by applying the actions recommended by the approach and the total number of the test cases for which the approach successfully proposed a solution (improving or degrading performance compared to the user's actions).

B.7 IMPLEMENTATION

- *Software environment*: All steps of the experimental process (the physical model of the office used for the generation of the database, the similar case retrieval process, and the adaptation process) are implemented using the Python 3.9 language running on Windows 10 Professional.
- *Computation material*: We have run both the dataset generation and the evaluation tests on a laptop equipped with an Intel[®] Core[™] i7-8559U CPU (having 2.70GHz and 2.60 GHz clock speed) and 16 GB of RAM.

C SIMILAR CASES RETRIEVAL

We use the approach given in Boulmaiz et al. (2021b) to estimate the similarity and define the similar source cases in the training set to each target case in the test set. The similarity consists in defining a threshold context distance σ beyond which the cases are not considered similar, in this case, the similar cases are all the source cases (learning cases) having a context distance to the target case (test case) below the threshold σ .

The authors propose a two-step approach to identify source cases similar to the target problem:

PARAMETER	VALUE
Population size	50
Chromosome size	12
Mutation rate	0.01
Crossover rate	0.8
Selection method	Roulette wheel
Generations	500

Table 5: Genetic algorithm parameters

1. a combined method of a genetic algorithm and a clustering algorithm is used to weight the variables.

The weighting process starts by initializing a vector with random values between 0 and 1 whose sum is 1, corresponding to context and actions variables' weights. The vector length corresponds to the number of context and action variables to be weighted (12). The context and action variables thus weighted are used to cluster the cases in the case base using Kmeans based on a context-action Euclidean distance as shown in Equation 3.

An optimization approach based on a genetic algorithm is used to find the optimal values of the weights. The optimization process (the fitness function) consists in minimizing both the average of the context-action distance and the average of the Euclidean distance between the effect variables of the cases in each cluster. Table 5 gives the parameters of the genetic algorithm.

For each test set, the genetic algorithm is applied to the validation set (the four other sets of the k-fold cross validation, see Section B.4.2).

2. a statistical method based on F1-measure is applied to determine the threshold of the context distance for which the source cases are considered similar to the target case. We rely on the context distance to determine the source cases (from learning set) similar to the target cases (from test set) using the context variables because the latter are the only data available for the target case.

The objective is to learn the maximum context distance between the source cases (learning set) for which two cases are considered similar. This distance is used as a threshold σ to determine source cases similar to the target case (from the test set) by computing the context distance between the latter and the source cases.

The rationale of this approach is to use the F1-measure to determine the context distance σ between the cases in the learning set that maximizes both precision (proportion of selected similar cases to the total number of selected cases) and recall (ratio of selected similar cases to the total number of similar cases).

Columns 6, 7, 8, and 9 of Table 3 report the results for the threshold context distance related to each test set as well as the corresponding recall, precision, and F1-score.

D FURTHER EXPERIMENTS: REAL DATASET

D.1 SIMILAR CASES RETRIEVAL

Recall that the real dataset validation process uses the LOOCV method (see Section B.4.1) where each case in the case base is used as a test set (target case) and the rest of the cases (97 cases) are used as a learning set (source cases). This is repeated such that each case is, in one iteration, the test case. Applying the approach described in Section C to extract source cases similar to the test case, generates as many context distance thresholds as the number of cases in the case base. This is illustrated in Figure 5, which depicts the context distance threshold for each iteration corresponding to a single test case. While the maximum threshold reaches the value of 1.208 for the test case with index 84, the minimum threshold distance is obtained for the case of index 51 with the value of 1.134. Figure 6 shows the F1-score, precision, and recall obtained for each threshold distance associated with each test case. The best score for the F1-measure is obtained for the test case 54 (80.43%) with a precision of 79.72% and a recall of 81.15%. The worst score is 58.41% which is recorded for test case 81 where the precision is 57.94% and the recall is 58.897%.



Figure 5: Context distance threshold for real dataset.



Figure 6: F1-measure, recall, and precision of each distance threshold for real datset evaluation.

Table 6 shows the evaluation results using the real dataset. The analysis of the results confirms the findings found in the analysis of the validation process on the synthetic dataset

Tuble 6. Summery of results on real autuset									
	TEST SET	S1 METRICS							
AITKOACII		PER (%)	APR(%)	EQR(%)					
CBR - S		19.73	50.23	50.23					
CBR - B		20.73	49.85	49.85					
CBR - P		26.12	39.64	48.09					
CBR - R		13.54	28.80	37.95					
CBR - APF		30.23	100	100					

 Table 6: Summery of results on real dataset