

# A decision support framework for robust R&D budget allocation using machine learning and optimization

Hoon Jang

Center of Public R&D Program Evaluation, Science and Technology Policy Institute, Sejong 30147, South Korea



## ARTICLE INFO

### Keywords:

Research and development  
Data-driven R&D budget allocation framework  
Public R&D program  
Machine learning  
Robust optimization

## ABSTRACT

Considering that government funding agencies make decisions on research and development (R&D) budget allocation to support an increasing number of research proposals, effective decision support systems are necessarily required. Motivated by the efforts of the Korean government, we propose a new decision support framework for allocating an R&D budget such that it maximizes the total expected R&D output. The proposed framework incorporates an R&D output prediction model with an optimization technique. We first employ a machine learning algorithm to accurately estimate future R&D output. Then, we apply a robust optimization technique to hedge against uncertainty in the predicted R&D output values. If not properly accounted for, this uncertainty may yield an inefficient budget allocation plan, thus hindering the operation of the R&D budgeting system. We demonstrate the effectiveness of the proposed model by applying it to a national R&D program conducted by the Korean government. Specifically, using the same budget, our budget allocation plan can achieve an output 13.6% greater than the actual R&D output. Thus, our model helps to attain *allocation efficiency* by systematically allocating budgets. We also observe the *price of robustness* when our model conservatively allocates budgets in order to hedge against uncertainty in the R&D predictions. Our findings offer insights for both policymakers and researchers related to designing better budget allocation systems for national R&D programs.

## 1. Introduction

Governments typically aid public R&D activities via funding agencies. Thus, such agencies spend a significant amount of researchers' R&D budgets. In a typical R&D funding procedure, agencies have to decide how to allocate R&D budgets in order to maximize the return on their investment. For instance, in the United States, the National Institutes of Health (NIH), the largest public R&D funding source worldwide for biomedical research, is continually trying to improve its funding mechanism to support the “right” number of projects with the “right” budgets. However, as noted by Park et al. [1], there is still room for improvement in allocating R&D budgets effectively.

In general, there are inherent difficulties in allocating R&D budgets optimally. For example, there is often a significant lag before the economic benefits of the R&D are realized. The exact amount of resources (e.g., time, # of researchers) required to handle an R&D project task is difficult to estimate. A success is uncertain, and the benefits to society of research findings may not be clear. Moreover, R&D budget allocation decisions are becoming more complex for decision-makers and R&D funding authorities.

Of the aforementioned difficulties, in this study, we focus on the

uncertainty related to the return on investment, which we refer to as R&D output. We do so because a key motivation of public R&D funding is to yield value through knowledge creation and innovation. Though important, we find that current funding practices do not consider this property appropriately [2–4]. For instance, the Korean government is experiencing insufficient returns on its public R&D investment. The amount allocated to R&D budgets by the government has increased significantly. In 1964, the total was \$1.79 million (USD) ( $\approx 2$  billion (KRW)), and is expected to surpass \$17.9 billion USD ( $\approx 20$  trillion (KRW)) for the first time in 2019. The Korean government's R&D budget is ranked fourth (in terms of aggregate amount), after those of the United States, Japan, and Germany. The annual average budget increase for R&D over the past 10 years (2008–2019) is approximately 6%, which is the fifth largest of the OECD member countries. However, the R&D output has not kept pace with the growth in the amount allocated to the budgets. For example, the numbers of outstanding publications and triadic patent families, which are considered key performance measures of R&D investment, have been stagnant over the past 10 years compared with the size of the government's R&D investment [5].

To overcome this problem, the Korean government has expended significant effort. For example, they established the Office of Science

E-mail address: [hoonjang@stepi.re.kr](mailto:hoonjang@stepi.re.kr).

<https://doi.org/10.1016/j.dss.2019.03.010>

Received 2 January 2019; Received in revised form 27 March 2019; Accepted 27 March 2019

Available online 05 April 2019

0167-9236/ © 2019 Elsevier B.V. All rights reserved.

and Technology Innovation in 2004 (reopened in 2017) to allocate and adjust the budgets of R&D programs carried out by Korean government. Furthermore, they now compile, allocate, and manage their R&D budget in a systematic and efficient manner. Recently, the government has also begun studying scientific R&D budget allocation techniques to more effectively manage their budget.

Motivated by the Korean government's ongoing efforts, we design a decision support framework for determining an optimal R&D budget allocation that maximizes the expected R&D output. Specifically, in the context of R&D, the uncertainty of future success demands that we develop rigorous decision support systems for allocating R&D budgets by precisely estimating each project's output. For this purpose, we devise a hybrid decision support framework that combines a machine learning algorithm with an optimization technique. Machine learning is used to estimate the expected R&D output of each project as accurately as possible. Then, we use a robust optimization to allocate R&D budgets optimally based on an expected deviation in the prediction values. By using the proposed framework, we conduct a case study to verify its applicability and validity. For this, we use national-level data from a well-known Korean R&D program.

The rest of this paper is structured as follows. In Section 2, we review the literature related to this study. In Section 3, we introduce the proposed framework for allocating R&D budgets to maximize the expected R&D output. In Section 4, we discuss our case study, data, and experiments. Section 5 describes the results of the experiments, and Section 6 discusses these results. Finally, Section 7 concludes the paper.

## 2. Literature review

This section reviews previous studies in fields pertinent to this study. Budget allocation (or project selection, if the scope is expanded) has been studied for a long time in fields such as finance and project management, in addition to R&D. Owing to the extensive volume of related studies, we focus on studies that allocate public R&D budgets using a mathematical approach. Specifically, we first review prior studies that predict R&D output (or outcomes) in Section 2.1. Next, we review existing budget allocation models in the context of public R&D (Section 2.2). Then, we briefly review robust optimizations and their applications to public R&D budget allocation (Section 2.3). In Section 2.4, based on our review of prior studies, we describe the novelty of our work. For a general overview of allocating R&D budgets (or project selection), refer to Refs. [6–8].

### 2.1. Studies on predicting R&D output

In general, it is fundamentally difficult to estimate future R&D output, largely because R&D is related to various uncertain factors. Moreover, it is difficult to determine a clear, causal relationship between such factors and R&D output. Thus, previous studies have attempted to use various methods, including expert consultations and index analyses, as qualitative methods [9–11]. These methods are particularly helpful when evaluating new R&D programs for which it is difficult to obtain data or that do not have sufficient data. However, owing to the nature of qualitative methods, they require a legitimate selection of experts and sufficient background knowledge. In contrast, quantitative methods such as regression analyses, time series analyses, and data envelopment analyses (DEA) can be used when there are sufficient data. However, these methods are not easy to use when the data cannot be expressed mathematically [12–15].

Recently, more diverse methods are being used to analyze the effects of R&D investments. These methods include the system dynamics simulation model, which considers the complicated interaction of various factors [16–18]. The latter method is useful for analyzing the causal relationships between key factors. However, it is difficult to establish the model when there are many variables, and it is not always guarantee reliability. More recently, deep neural networks (DNNs) have

been used to analyze the effects of R&D investments [19–24]. Compared with other quantitative methods, DNNs show excellent performance when complex correlations exist between variables. DNNs are also used to estimate future R&D output or risks at an early stage. However, most studies that use artificial neural networks apply to the construction industry, and thus are not widely applicable [20–22]. Moreover, to the best of author's knowledge, no studies focus on making rigorous R&D budget allocation decisions by using artificial neural networks.

### 2.2. Studies on R&D budget allocation

Studies on R&D budget allocation models date back to the 1950s. Lorie and Savage [25], Markowitz [26], and Bernhard [27] conducted seminal works on allocating assets by evaluating a set of projects using a simple financial evaluation index. As optimization theory and computing power have improved, various mathematical approaches to R&D budget allocation problems have been proposed. Heidenberger and Stummer [28] classified R&D budget allocation and project selection models into six descriptive methods (mathematical programming, simulation based optimization, heuristics, real-option, etc.). They then classify previous studies on mathematical programming based on the characteristics of their mathematical programming (i.e., linear, nonlinear, multi-period, or probabilistic).

Ghasemzadeh and Archer [29] proposed a decision-making support system for R&D budget allocations based on an integer linear programming and analytic hierarchy process (AHP). Badri et al. [30] presented an R&D project selection model for medical care that uses goal programming under a constraint on the total budget amount. This study is valuable because it mathematically solves the problem by considering both the utility and cost of research and the project implementation period. Coldrick et al. [31] proposed a method based on traditional portfolio methods (scoring methods, such as discounted cash flow, based on real-option theory), which they applied to R&D project selection. Gabriel et al. [32] presented an optimal project selection model that considers the possibility of success of an R&D project. Their proposed model combines a multi-objective optimization model, Monte Carlo simulation model, and AHP. Doerner et al. [33], Abdelaziz et al. [34], Medaglia et al. [35], Tolga [36], and Carazo et al. [37] also studied multi-objective optimization models in R&D budget allocation and project selection problems.

More recently, studies have begun considering the complicated characteristics inherent in R&D. Guo et al. [38] proposed using 0–1 nonlinear mathematical programming to select R&D projects, considering the four attributes of R&D (output, resources, technology, and risk). Gutjahr et al. [39] studied an R&D project selection model in which they consider the possibility of success of a project and include a penalty for exceeding the project implementation period. Litvinchev et al. [40] explored an R&D project selection model for public organizations. They developed an optimal R&D project selection model that considers the balance between the qualitative elements of a project and the total number of projects selected. Jung and Seo [41] presented an analytic network process approach for evaluating R&D projects with heterogeneous types of objectives. Aryanezhad et al. [42] and Bhattacharyya et al. [43] proposed portfolio selection models that use a fuzzy optimization method. Luo [44] developed an optimal R&D project selection model that controls for potential risk factors related to the market and technology development in the implementation stage of a project. Casault et al. [45] analyzed various R&D project selection methods, where they considered the interactions among R&D projects and proposed an R&D project selection method by considering both quantitative and qualitative evaluation measures equally. Costantino et al. [46] conducted a seminal work using an artificial neural network to assess the future output of an R&D project. However, it is worth noting that they did not use mathematical models to select the best projects. Arratia-Martínez et al. [47] proposed a mathematical model

for selecting R&D projects that considers four characteristics (social objective, emphasis areas, geographical influences, and nonmonetary factors) in the domain of public organizations.

### 2.3. Brief review of robust optimization

As covered in many studies, uncertainty must be handled appropriately in R&D budget allocation and project selection problems. Recently, stochastic programming and robust optimization models have been applied in related research fields. The robust optimization model and stochastic programming provide a popular method of determining an optimal solution in problems involving uncertainty. In particular, this method can be used effectively when the data are random or when it is difficult to estimate the mathematical distribution of the uncertain variables owing to environmental changes or a lack of knowledge about the variables. In other words, even when it is difficult to know the probabilistic distribution of the data used in the optimization, the robust optimization method can be applied successfully. This method was first proposed by Soyster [48]. Since then, various models and solutions have been presented. Note that robust optimization has received much attention, from both academic and industrial areas, as a result of theoretical and computational development since the 2000s.

Among the many approaches in robust optimization, a model proposed by Ben-Tal and Nemirovski [49] is widely used. They proposed a mathematical formulation using the conic quadratic robust optimization model, such that the ellipsoidal uncertainty can be considered appropriately in the mathematical programming. The  $\Gamma$ -approach proposed by Bertsimas and Sim [50] is another popular robust optimization model. By assuming that there is almost no chance that all variables fall into the worst case, they proposed a mathematical formulation allowing users to control the robustness of the model. Refer to Melo et al. [51], Gülpinar and Pachamanova [52], and Govindan et al. [53] for further details.

Relatively few studies apply the robust optimization model to R&D budget allocation or project selection models. Hassanzadeh et al. [54] developed a multi-objective optimization model for R&D project selection in which they assume that uncertainty is inherent in the variables to be included in the constraints and the objective function of the optimization model. Then, they solved this problem using the robust optimization model. Around the same time, Mild et al. [55], Bekiros et al. [56], and Liu et al. [57] developed budget allocation and project selection models using the robust optimization model. In particular, Liu et al. [57] considered the cognitive bias of the project selection decision maker, and proposed a solution applying the particle swarm optimization method, considering the complexity of the model. Note that, in general, the budget allocation models that employ a robust optimization technique assume input parameters based on simple statistics from historical data [54,55,58].

### 2.4. Contribution of our work

Based on the above review, several types of studies in fields such as R&D output estimation, R&D budget allocation, and robust optimization are related to our study. However, although the aforementioned areas have been studied thoroughly, we believe that our work will enrich each of them in two ways. First, although many studies have been conducted on estimating the output of a project or on project selection, few studies propose a combination of a prediction model and an optimization method. Specifically, we use a machine learning algorithm to estimate the expected R&D output of each project. Then, we apply the resulting values to optimally allocate R&D budgets. In this step, we use a robust optimization to hedge against uncertainty in the predicted values. We can thus enhance the effectiveness of our approach (combining the prediction model with the optimization method), because robust optimization can properly address the uncertainty of a prediction model. To the best of our knowledge, we are

the first to propose such an integrated framework for designing an optimal budget allocation plan comprising public R&D projects. Second, we conduct experiments to examine whether a data-driven budget allocation scheme can be used successfully in an actual R&D program. The data are provided by the Korean government. The results verify the applicability and validity of the proposed framework, which we believe that we can aid policymakers and researchers in expanding the scope of their work.

## 3. Proposed modeling framework

In this section, we propose a decision support framework for optimally allocating an R&D budget. In Section 3.1, we introduce a prediction model for estimating the expected R&D output. In Section 3.2, we propose an R&D budget allocation model that maximizes the sum of the expected R&D outputs. By assuming variability exists in the estimated R&D output values, we build a mathematical model using a robust optimization technique.

### 3.1. R&D output estimation model

In order to estimate the expected R&D output of each project, we use a machine learning algorithm that has recently shown remarkable performance. Specifically, we use AutoML, an abbreviation for automated machine learning, a machine learning algorithm in which the machine finds the solution to the given data and produces results on its own. AutoML is a user-friendly machine learning platform. Unlike traditional machine learning algorithms, AutoML eliminates or minimizes most of the parts that previously required a user's involvement (e.g., adjusting parameters). AutoML automatically tunes all related parameters based on embedded algorithms. More importantly, AutoML performs sufficient well, but with significantly less user involvement.

We use AutoML in H2O, an open source, in-memory, distributed, fast, and scalable machine learning and predictive analytics framework. At a general level, we prepare a training dataset and define a small number of parameters (e.g., maximum calculation time, performance evaluation criteria). Then, AutoML provides the best model based on various machine learning algorithms. Specifically, we apply five algorithms: distributed random forest, generalized linear model, gradient boosting machine, deep learning, and stacked ensemble.

### 3.2. R&D budget allocation model

This section presents a mathematical model for the optimal allocation of R&D budgets according to expected R&D outputs. The purpose of the mathematical model is to maximize the sum of the expected R&D outputs of the projects. We use the following notation:  $i \in I = \{1, 2, 3, \dots, n\}$  denotes a project; and  $j \in J = \{1, 2, 3, \dots, m\}$  represents a budgeting option that can be allocated to a project. In other words, we allow each project  $i$  to receive different budget allocations, and the final allocation is determined by the R&D output value predicted by the R&D output estimation model. For instance, if project  $i$  is expected to perform well, it will receive an increased budget allocation compared to its original budget. On the contrary, if project  $i$  is expected to perform worse than others, it will be allocated a reduced budget compared to the original amount. The parameters included in this problem are  $e_{ij}$ ,  $b_{ij}$ , and  $B$ , where  $e_{ij}$  represents the expected R&D output when project  $i$  is given the  $j$ th budgeting option. Similarly,  $b_{ij}$  represents the budget amount when project  $i$  is assigned the  $j$ th budgeting option. Lastly,  $B$  is defined as the total budget amount that can be used by the R&D program. In other words,  $B$  is the total available budget that can be assigned to the R&D program. We define the principal decision variable,  $y_{ij}$ , as a binary variable, taking the value one when the  $j$ th budgeting option is assigned to project  $i$ , and zero otherwise.

Using the above decision variable and parameters, we construct an R&D budget allocation model as follows:

$$(P^O) \max \sum_{i \in I} \sum_{j \in J} e_{ij} y_{ij} \quad (1)$$

$$\text{subject to } \sum_{i \in I} \sum_{j \in J} b_{ij} y_{ij} \leq B \quad (2)$$

$$\sum_{j \in J} y_{ij} = 1 \quad \forall i \in I \quad (3)$$

$$y_{ij} \in \{0, 1\} \quad \forall i \in I, \forall j \in J \quad (4)$$

The objective function of the model ( $P^O$ ) maximizes the total expected R&D output from  $n$  projects. The model has three constraints. First, the upper limit of the total budget for the program is limited to  $B$  (constraint (2)). Each project  $i$  must be assigned one of  $m$  budgeting options. In other words, projects that are not selected are not considered in this model (constraint (3)). The decision-making variable must be binary (constraint (4)).

With this nominal model, we propose a robust version of the R&D budget allocation model, assuming uncertainty in the expected R&D output predicted by AutoML. In other words, assuming there is uncertainty inherent in the predicted values of R&D output, we propose a robust budget allocation model that avoids the negative consequences of such uncertainty. Specifically, we use the  $\Gamma$ -approach proposed by Bertsimas and Sim (2004). As noted above, the  $\Gamma$ -approach enables us to control the robustness of the model, allowing us to avoid arriving at too conservative a solution.

We assume that the uncertainty of the expected output for each R&D project follows a box-type uncertainty set. Such a set assumes that the parameters designated by users can be revealed randomly between the average ( $\bar{e}_{ij}$ ) and deviation of the original value ( $\hat{e}_{ij}$ ):

$$A = \{e_{ij} \in E \mid e_{ij} \in [\bar{e}_{ij} - \hat{e}_{ij}, \bar{e}_{ij} + \hat{e}_{ij}]\}. \quad (5)$$

Note that box-type uncertainty sets are the simplest of the polyhedral uncertainty sets, but are widely used in robust optimization. In particular, they provide useful approximations when little is known about the uncertain parameters, because they provide the greatest scope of possibility when realizing the parameters. Thus, we decide use box-type uncertainty sets to analyze uncertainty in the expected R&D outputs, because we have little information on the expected output of each project. Under this uncertainty set, the robust model ( $P^R$ ) for the optimal allocation of R&D budgets is given as follows:

$$(P^R) \max \sum_{i \in I} \sum_{j \in J} \bar{e}_{ij} y_{ij} + \Delta(y, \Gamma) \quad (6)$$

$$\text{subject to } \sum_{i \in I} \sum_{j \in J} b_{ij} y_{ij} \leq B \quad (7)$$

$$\sum_{j \in J} y_{ij} = 1 \quad \forall i \in I \quad (8)$$

$$y_{ij} \in \{0, 1\} \quad \forall i \in I, \forall j \in J. \quad (9)$$

As shown above, this model is not complicated. However, it includes an inner optimization problem,  $\Delta(y, \Gamma)$ , called a protection function, which helps to consider the worst-case scenario involving the uncertainty in the output. Specifically, we use  $s_{ij}$  to control the worst (i.e., largest) deviation that can be considered by the protection function. We first formulate the protection function on the left-hand side, as follows:

$$\begin{aligned} \Delta(y, \Gamma) = \min & \sum_{i \in I} \sum_{j \in J} \hat{e}_{ij} y_{ij} s_{ij} \\ \text{subject to} & \sum_{i \in I} \sum_{j \in J} |s_{ij}| \leq \Gamma \\ & |s_{ij}| \leq 1, \quad \forall i, \forall j \\ \Leftrightarrow & -\max \sum_{i \in I} \sum_{j \in J} \hat{e}_{ij} y_{ij} s_{ij} \\ \text{subject to} & \sum_{i \in I} \sum_{j \in J} s_{ij} \leq \Gamma \\ & 0 \leq s_{ij} \leq 1, \quad \forall i, \forall j. \end{aligned}$$

Because the left-side of the protection function is the minimization problem, it is evident that optimal values of  $\Delta(y, \Gamma)$  will have  $\Gamma$  of  $s_{ij}$  at  $-1$ , and all other  $s_{ij}$  at zero. Therefore, we can equivalently rewrite the left-hand side as shown on the right-hand side, without the absolute function for  $s_{ij}$ . Note that the rewritten formulation (the right-hand side) is linear, feasible, and bounded for a fixed value of  $y$ . Therefore, we can now demonstrate that the robust model ( $P^R$ ) has an equivalent linearized model ( $P_L^R$ ). For further details, refer to Proposition 1.

**Proposition 1.** ( $P^R$ ) has an equivalent linear formulation, ( $P_L^R$ ).

$$(P_L^R) \max \sum_{i \in I} \sum_{j \in J} \hat{e}_{ij} y_{ij} - q\Gamma - \sum_{i \in I} \sum_{j \in J} r_{ij} \quad (10)$$

$$\text{subject to } \sum_{i \in I} \sum_{j \in J} b_{ij} y_{ij} \leq B \quad (11)$$

$$\sum_{j \in J} y_{ij} = 1 \quad \forall i \in I = 1 \quad \forall i \in I \quad (12)$$

$$q + r_{ij} \geq \hat{e}_{ij} y_{ij} \quad \forall i \in I, \forall j \in J \quad (13)$$

$$q, r_{ij} \geq 0, y_{ij} \in \{0, 1\} \quad (14)$$

**Proof.** First, it is true that the inner optimization problem  $\Delta(y, \Gamma)$  is feasible and bounded for fixed values of  $y$ . Therefore, we now consider the dual of the inner optimization problem:

$$\begin{aligned} \Delta(y, \Gamma) = -\max & \sum_{i \in I} \sum_{j \in J} \hat{e}_{ij} y_{ij} s_{ij} \\ \text{subject to} & \sum_{i \in I} \sum_{j \in J} s_{ij} \leq \Gamma \\ & s_{ij} \leq 1, \quad \forall i, \forall j \\ & -\min q\Gamma + \sum_{i \in I} \sum_{j \in J} r_{ij} \\ \Leftrightarrow & \text{subject to } q + r_{ij} \geq y_{ij} \quad \forall i, \forall j \\ & q, r_{ij} \geq 0, \quad \forall i, \forall j. \end{aligned}$$

By strong duality, the dual problem of  $\Delta(y, \Gamma)$  is also feasible and bounded, and the optimal value should be the same as that of the primal model. Using this property, we now replace ( $P^R$ ) with ( $P^R$ )

$$(P^R) \max \sum_{i \in I} \sum_{j \in J} \bar{e}_{ij} y_{ij} - \min_{\substack{q+r_{ij} \geq \hat{e}_{ij} y_{ij}; (15) \\ q, r_{ij} \geq 0}} (q\Gamma + \sum_{i \in I, j \in J} r_{ij})$$

subject to Eqs. (7)–(9), which is equivalent to ( $P_L^R$ ). To demonstrate this, first suppose that we solve ( $P_L^R$ ) and obtain an optimal solution. Because ( $P_L^R$ ) is solved optimally, the optimal solution is also feasible for ( $P^R$ ), and the optimal value will be same. Conversely, suppose we have an optimal solution to ( $P^R$ ). Because it is optimal for ( $P^R$ ), it is also valid for the constraints shown in the inner optimization problems. That is, an optimal solution for ( $P^R$ ) is feasible for ( $P_L^R$ ) and gives the same optimal value. Therefore, we can discard the “min” term from ( $P^R$ ), which makes it equivalent to ( $P_L^R$ ).  $\square$ .

Note that ( $P_L^R$ ) is now formulated as an integer program; therefore, we assume that it can be handled by using a commercial solver, such as CPLEX.

#### 4. Case study

In this section, we apply the proposed framework to a real-world case and verify its applicability. In Section 4.1, we give a general overview of the case study, and in Section 4.2, we present our data. In Section 4.3, we describe the design of our experiments.

##### 4.1. Problem description

The purpose of this case study is to allocate a budget to each project within an R&D program to maximize the overall expected R&D output by applying the proposed framework. In doing so, we verify the applicability and utility of the framework. The program specified in this



study is BrainKorea 21 Plus (BK21 Plus), organized by the Ministry of Education, which is a follow-up program to BrainKorea 21 (BK21). The previous program was one of many national programs aimed at nurturing world-class graduate schools and human resources. In all, approximately \$3.1 billion (USD) (= 3.5 trillion (KRW)) was invested over 14 years, from 1999 to 2012. BK21 Plus represents Phase II of BK21 (planned for a total of seven years, from 2013 to 2020), aiming to “nurture master and doctorate-level creative manpower to promote future national competitiveness and support creation of new knowledge and technology based on creativity.”

The key performance measures of the BK21 Plus program are the number of publications, number of patent registrations, number of technology transfers and commercialization support for technology development output, nurturing of outstanding masters- and doctorate-level human resources, and employment rate for educational output. These are defined clearly and are easy to collect and measure in comparison with those of R&D programs conducted by the private sector. This is one of reasons why we chose this program for the case study.

#### 4.2. Data

To be realistic, it is necessary to use both qualitative data (e.g., expert consultations) and quantitative data, collected over many years in a rigorous manner. However, this is beyond the scope of our study. Instead, we use data collected from the National Science & Technology Information Service (NTIS). The NTIS is currently operating the national R&D information standard database that embraces all government departments. It also collects, processes, and discloses 422 information standards for the major special agencies (17) and project management agencies (125) that handle all national R&D information for each department. There are 628,350 cases of data related to R&D projects, and 4532 cases of R&D output provided by the NTIS as of October 2017.

We use the data of the BK21 Plus program for the period 2013 to 2016. Before 2016, the NTIS discloses outputs as the number of publications and the number of patents only. Thus, it is not possible to estimate other types of outputs (e.g., technology transfer). We also consider the gap between when the R&D was conducted and when publications (and patents) were published. For simplicity, we set the maximum allowable gap to three years. This means that we focus on R&D projects that proceeded in 2014, and use three-year data to analyze the R&D outputs of the projects conducted in 2014. The variables used in this case study and their descriptive statistics are shown in Table 1.

#### 4.3. Experimental setting

We design the experiments in two phases. In the first phase, we determine an optimal budget allocation plan without considering the uncertainty in the expected R&D output. Here, we explore many cases by varying the total budget,  $B$ . Specifically, we change the total budget from  $-20\%$  to  $+30\%$  of the original volume reported by the data, in increments of  $5\%$ . For each scenario, we examine the performance of our framework by measuring  $\rho$ . In the second phase, we analyze the impact of using the proposed robust budget allocation model. Here, we obtain optimal budget allocation plans from the robust model ( $P_L^R$ ) and the nominal model ( $P^O$ ) by varying the range of deviation in the expected R&D output ( $\hat{e}_{ij}$ ) and the level of robustness ( $\Gamma$ ). In total, we prepare 42 different scenarios. In this phase, to verify the effectiveness of the robust model, we conduct a simulation study with 100 randomly generated cases, varying the values of the R&D outputs.

The primary indicator of our experiments is defined as the rate of improvement in the expected R&D output compared with the original R&D output calculated from the actual data:

$$\rho = \frac{z^* - z}{z}, \quad (15)$$

where  $z$  is the R&D output from the real data, and  $z^*$  is the objective function value from the proposed framework. Note that we consider seven budgeting options for each R&D project (cardinality of  $J$  is equal to seven). Specifically, considering a practical constraint that the budget of each project cannot be changed significantly from its original level, we set the range as  $\pm 30\%$  of its original volume. The basic unit is  $10\%$ , including  $0\%$ , which represents maintaining the original budget.

As mentioned previously, we use the H2O platform to estimate the R&D output, and use the open source statistical program R. For the budget allocation, the optimization program CPLEX v.12.8 is used, and the computer programming language Java is used for operation. The computing environment used in the experiment is as follows: Intel i7-6700 (3.4 GHz), 8 GB RAM, Windows 7 Professional. We arbitrarily set the maximum calculation time of AutoML to 3 h. This is warranted by the results of the preliminary tests to determine the maximum allowable time for AutoML. The results show no significant changes in the quality of prediction results after 3 h. We also set the maximum calculation time for the budget allocation model to an hour. Though we set the limit for this calculation, all instances are solved within a few minutes. In order to examine the computational burden of the proposed budget allocation model, we virtually increase problem sizes up to 300% of the original. The results show no computational burden (details are presented in Appendix A.1). However, it should be noted that since our model is formulated by binary programming, a computational burden may occur in some instances with large problem sizes. To handle such a difficulty, efficient solution algorithms, including meta-heuristic algorithms, can be employed to obtain a feasible solution of the problem.

### 5. Experimental results

Hereafter, we report the experimental results. In Section 5.1, we report the prediction results of the expected R&D output. Next, we present the R&D budget allocation results of the studied R&D program, with and without considering the uncertainty in the estimated values of R&D output (Sections 5.2 and 5.3, respectively).

#### 5.1. Results of estimating future R&D output

We describe the prediction results of the expected R&D output in Table 2. In this study, we primarily use root mean square error (RMSE) which is commonly employed to examine the quality of prediction models. Additionally we report the mean absolute error (MAE). Note that the values in this table are five-fold cross-validation values. As shown in Table 2, five machine learning algorithms are used to predict the expected R&D output. Of these, AutoML is superior to other benchmarking algorithms, showing satisfactory performance. It is worth noting that though we use RMSE to construct a model to predict future R&D output, MAE also consistently indicates that AutoML outperforms other benchmarking algorithms.

Specifically, the values of RMSE and MAE using AutoML are 39.90 and 20.80, respectively, or 48.1% and 120.5%, respectively, better than other methods. Moreover, using a validation data set, we found a high correlation between the actual R&D output and the estimated R&D output using the model generated by AutoML ( $R^2 = 0.869$ ). Therefore, it is possible to conclude that accepting our approach using AutoML to evaluate the expected R&D output is warranted.

#### 5.2. Results of allocating R&D budgets without considering uncertainty

We next illustrate the effectiveness of using our model when allocating R&D budgets. We first present the level of improvement in expected R&D output that can be attained when allocating R&D budgets using our proposed framework (Fig. 1). In this experiment, we arbitrarily vary the total budget amount that can be used for the BK21 Plus R&D program. Note that the horizontal axis indicates the increment rate

**Table 1**  
Variables applied to AutoML and the descriptive statistics of variables ( $n = 458$ ).

	Variable	Type	Values	Percentage (%)
Input	Research Organization	University	458	100%
		Research institute	0	0%
		Private company	0	0%
		Others	0	0%
	PI's Major Research Area	Engineering	202	44.1%
		Science	130	28.4%
		Medicine	48	10.5%
		Others	78	17.0%
	Research Phase	Basic	186	40.6%
		Applied	200	43.7%
		Development	16	3.5%
		Others	56	12.2%
	NCCST <sup>a</sup> (Level I) (in total 33 codes)	EE (Information & Telecommunication)	44	9.6%
		NC (Chemistry)	44	9.6%
		EA (Mechanical Engineering)	43	9.4%
		Others	327	71.4%
	NCCST (Level II) (in total 371 codes)	EA02 (Manufacturing based technology)	21	4.6%
		NA09 (Applied statistics)	17	3.7%
		EE11 (Telecommunication/information technology)	16	3.5%
		Others	404	88.2%
	Socio-economic Objective Code (in total 13 codes)	X12 (Progress in knowledge)	266	58.1%
		X01 (Health)	36	7.9%
		X02 (Education)	25	5.5%
		Others	131	28.5%
	Emerging Technology Area Code (Level I) (in total 6 codes)	Bio Technology	131	28.6%
		Information Technology	86	18.8%
		Nano Technology	57	12.4%
		Others	184	40.2%
	Emerging Technology Area Code (Level II) (in total 23 codes)	Healthcare application	57	12.5%
		Basic/Fundamental techniques	44	9.6%
		Information systems	33	7.2%
		Others	324	70.7%
	Total Amount of Fund (in 100 million (KRW))	–	2.67	
	Proportion of Labor (including stipend) Costs	–	45.7%	
	Proportion of Research Activities (e.g., Research Equipment, Materials, etc.) Cost	–	17.2%	
Output	Number of publications	–	42.4	
	Number of patents	–	28.1	

<sup>a</sup> National Classification Code of Science and Technology.

**Table 2**  
Performance ranking from AutoML against other benchmarking algorithms.

Algorithm	RMSE	MAE
AutoML (stacked ensemble)	39.90	20.80
Distributed random forest	43.32	27.15
Gradient boosting machine	50.86	28.77
Deep learning	46.21	59.76
Generalized linear model	80.64	50.76

of the total available budget from its baseline (0%). As noted, 0% refers to maintaining the current budget.

As shown in Fig. 1, it is theoretically possible to produce 13.6% more R&D output from the BK21 Plus program by re-allocating R&D budgets (keeping the current budget constant) using the proposed framework, as compared with the actual R&D output reported by the BK21 Plus program. This confirms that our framework achieves *allocation efficiency* in this R&D program. We also reduce the total amount of budget to 85% of its original volume, finding that the bar is still above zero, which means that we can achieve an improvement in the expected R&D output compared with the actual R&D output, even though 85% of the original budget is utilized. We also predict that, as the total budget increases, so will the expected R&D output<sup>1</sup>. However, the degree of

improvement is reduced as  $B$  increases. Overall, the above experimental results imply the importance of budget allocation efficiency. Using the same budget, our framework offers a way of increasing the productivity of R&D programs and avoiding projects a priori that can be predicted as poor.

To examine the results shown in Fig. 1 in more detail, we examine how much the budget changes in the projects. For this, we analyze five selective cases, including the baseline scenario, among 11 cases. For each case, we count how many projects' budgets are adjusted. The results are shown in Fig. 2. Overall, our proposed framework provides a budget allocation plan that gives more budget to projects that anticipate high expected R&D output, and reduces the budgets of projects that do not.

In the case of the baseline scenario, only 49 projects ( $\approx 10\%$  of the total number of projects) maintain their pre-assigned budget. In other words, 90% of the total projects receive a new budget, depending on their expected R&D output. Specifically, 182 projects (39.7% of the total number of projects) reduce their maximum permissible level ( $-30\%$  of its original budget, bars filled with black color) because the expected R&D output is poor. In contrast, 24% of the projects (111 of 458 projects) are assigned new budgets at the maximum allowable level ( $+30\%$  of its original budget, bars filled with hatched patterns).

Comparing the budget allocation plans of five problem instances, we find that the greater the available budget, the larger the number of projects allocated a higher budget (see the proportion of “add 10%”, “add 20%”, and “add 30%” in Fig. 2). To explain such a phenomenon, we first conjecture that the R&D budget of each project and the productivity (R&D output) has a positive correlation. Therefore, unless a

<sup>1</sup> Note that because we only have a point value for a comparison (when  $B = 0\%$ ), bars with hatched pattern are regarded as a virtual gain by the increase in the R&D budget. To more precisely examine the effect of our framework when  $B > 0\%$ , actual data on each level of  $B$  are required.

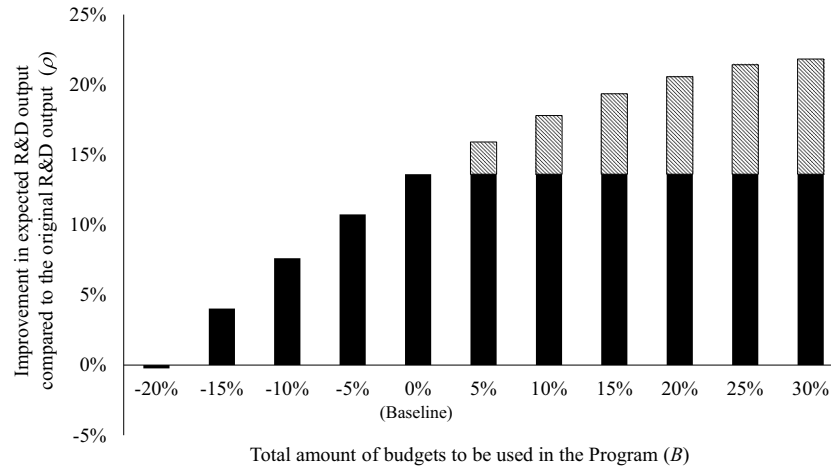


Fig. 1. Improvement in expected R&D output when the total budget can be varied from  $-20\%$  to  $+30\%$  by  $5\%$  each.

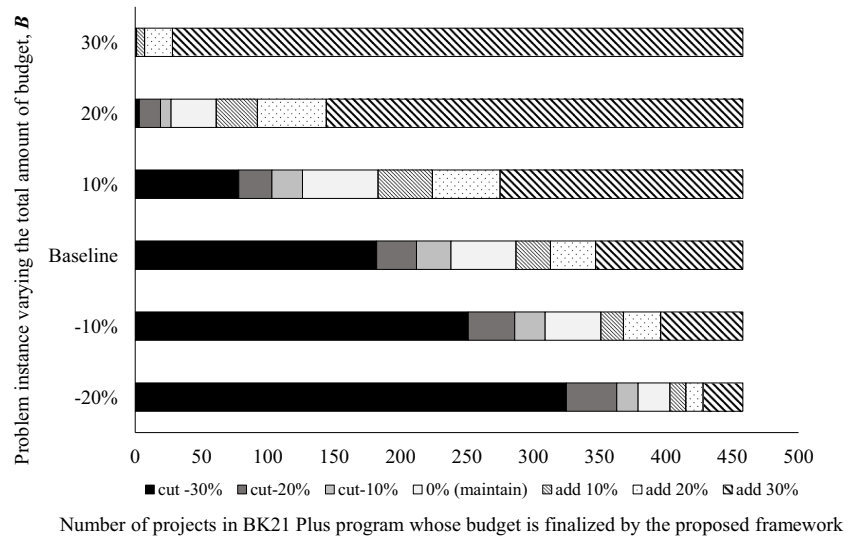


Fig. 2. Budget allocation plans of five selective problem instances.

project's productivity is expected to be reduced as the level of budget for the project is increased, it is natural that the number of projects that receive a greater budget will increase as the total available budget amount is raised. In contrast, if the total available budget is decreased under the baseline, the proposed framework finds a way of allocating budgets as effectively as possible. To this end, our proposed framework provides a budget allocation plan that cuts back most projects' budgets and concentrates on a few of projects that are expected to perform significantly better than others. In doing so, the proposed model tries to maximize the overall performance of the R&D program.

We now turn our attention to the sensitivity analysis. Specifically, we use the total sum of publications and patents as our primary performance measure of R&D output. However, it is plausibly assumed that each can have different weights. Therefore, we explore how changes of weights in the number of publications and patents affect the performance of the budget allocation solution. To examine this, we prepare five scenarios, varying the ratio of publication to patents from (1:3) to (3:1) in a baseline scenario, and evaluate the degree of improvement in the expected R&D output. The results confirm that there are no significant changes in the directionality, which means that although we give different weights to the number of publications and patents, at least some level of improvement compared with the actual R&D output is warranted. Details on the sensitivity results are depicted in [Appendix A.2](#).

### 5.3. Results of allocating R&D budgets, considering uncertainty

We now present the optimal budget allocation plan considering uncertainty of the expected R&D output. For this, we assume that the deviation of expected R&D output changes by 10% within the 10–30% range (Fig. 3a). The case in which the value of  $\Gamma$  that can be adjusted is also changed by 10% within a 10–30% range (Fig. 3b). As previously explained, the value of  $\Gamma$  is the parameter that adjusts the robustness of the solutions. If this value is zero, it does not consider the robustness of the solution at all, and thus is equivalent to the solution obtained from a nominal model ( $P^0$ ). If the value is one, the robust model considers the robustness as much as possible, and produces a solution that can be applied to all cases.

When the uncertainty of the expected R&D output is considered, the objective value decreases compared with that of the nominal model. The loss in the expected R&D output measured by the difference in objective value by the robust solution ( $z_{rob}$ ) relative to the nominal solution ( $z_{nom}$ ) ( $L = \frac{z_{rob} - z_{nom}}{z_{nom}}$ ) is shown in Fig. 3, and is even more significant if the deviation and value of  $\Gamma$  increase. Specifically, compared with the case in which the total budget is maintained, the objective value from the robust model is expected to decrease by approximately 5% when the deviation is assumed as 10%, compared with when uncertainty is not considered. Of course, when the deviation is assumed as 30%, the loss is computed as approximately 15%. This is

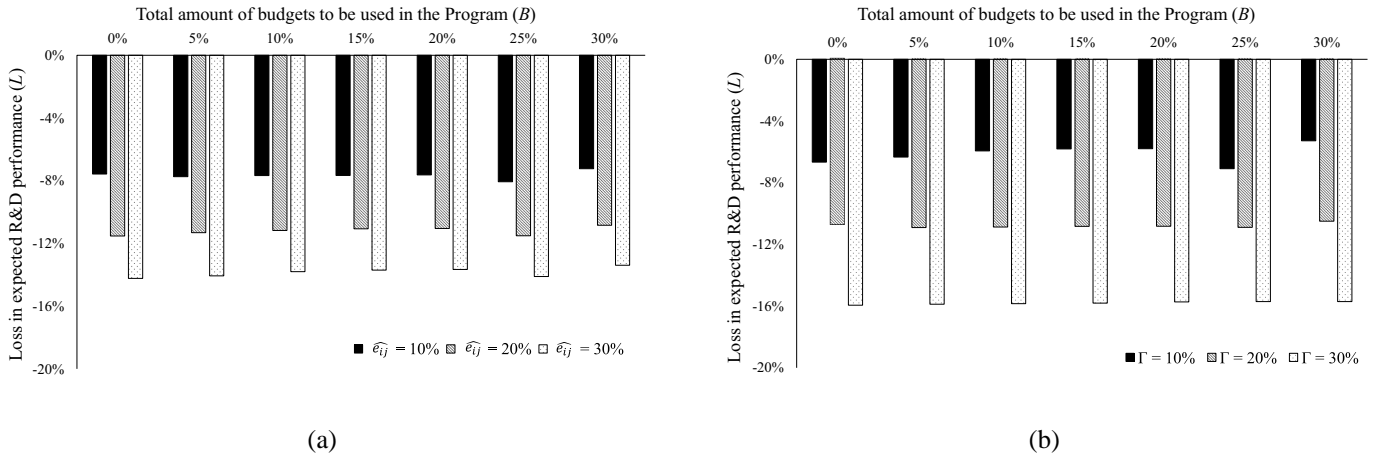


Fig. 3. Expected loss in R&D output ( $L$ ) when the uncertainty in the expected R&D output is introduced: varying the deviation of the expected R&D output ( $\hat{\epsilon}_{ij}$ ) (a); varying the value of  $\Gamma$  (b).

typically known as the *price of robustness*. Specifically, the results shown above can be explained by the characteristics of the solutions from the robust model. The robust model considers cases in which the expected R&D output turns out to be the worst in each project, and allocates the budget to maximize the expected R&D output accordingly. Thus, the budget is bound to be allocated very conservatively compared with when the uncertainty of expected output is not considered. However, the solutions obtained through the robust optimization model have high practical applicability in that they guarantee excellent performance when the uncertainty in the expected R&D output is assumed.

A closer examination provides a glimpse of how the solutions from the robust model work in practice, compared with those of the nominal model. For this, we conduct simulation experiments. Fig. 4 indicates the simulation results of three selective cases, controlling the level of  $\Gamma$  and the level of  $\hat{\epsilon}_{ij}$ . For instance, (Low, Low) in x-axis means that we set the lowest value of  $\Gamma$  ( $= 10\%$ ) and the lowest value of  $\hat{\epsilon}_{ij}$  ( $= 10\%$ ). With this assumption and the parameter setting, we measure the degree of improvement of budget allocation plans, with and without considering the uncertainty in the expected R&D output compared with the actual R&D output. As seen in Fig. 4, we report the performance as box-plots of each of two-budget allocation plan.

As shown in Fig. 4, if there exist deviations in the expected R&D

output of each R&D project, the budgeting plans obtained through the robust model generally outperform those from the nominal model, which uses the mean of R&D output. In other words, the results in Fig. 4 show the necessity of the robust model, in practice.

Specifically, Table 3 analyzes how the budget allocation plan obtained through the robust model is different from the budget plan obtained through the nominal model. The values in Table 3 show the number of projects, and the numbers in parentheses indicate the increase or decrease in budgets in each position. Specifically, with the nominal model, 182 projects are cut back their budgets to its maximum permissible level ( $-30\%$ ); however, after applying the robust model, 28 projects (among 182 projects) are assigned new budgets that are different from the maximum permissible level.

As shown in Table 3, there are cases in which the budget for each project increases or decreases when uncertainty is considered. Specifically, the 78 projects (above the diagonal line) receive an increased budget compared with the one when uncertainty is not considered; 61 projects (below the diagonal line) are given a decreased budget.

To examine this in more detail, we classify the projects with budgets that changed after applying the robust model into two groups: projects that increased the budget (group A); projects that decreased (group B). For each group, we calculate the average of the expected R&D outputs

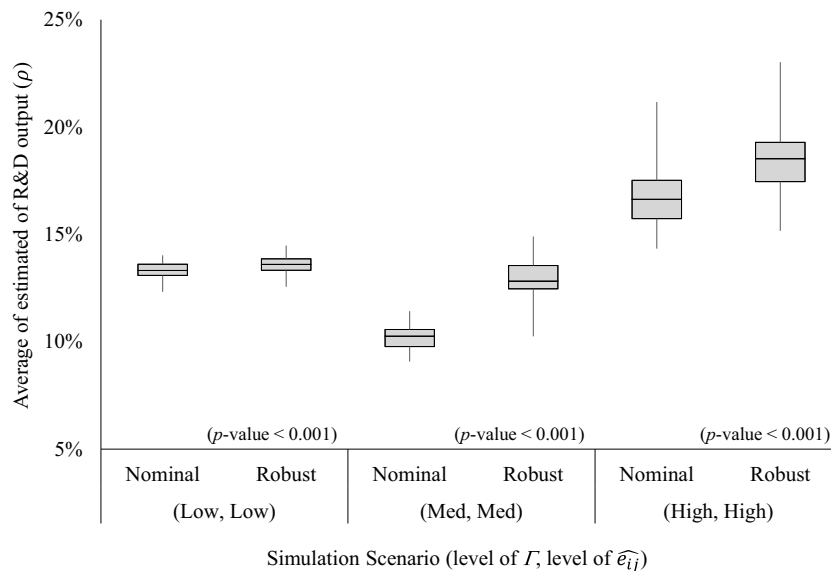


Fig. 4. Average of estimated R&D output when the uncertainty in the R&D output is realized.



**Table 3**  
Number of projects for which budgets have been changed: in case of (High, High) scenario.

		# of projects with changed budget when robust model is applied							Total
		–30%	–20%	–10%	0%	+10%	+20%	+30%	
# of projects with changed budget (with nominal model)	–30%	154	18 (4.27)	3 (3.28)	1 (1.51)	1 (0.42)	1 (0.54)	4 (6.01)	182
	–20%	2 (–0.1)	22	1 (0.35)	2 (3.67)	0	1 (2.54)	2 (3.17)	30
	–10%	1 (–0.72)	2 (–0.59)	11	7 (4.61)	3 (2.08)	0	2 (4.53)	26
	0%	1 (–2.35)	3 (–3.39)	4 (–1.84)	19	11 (5.91)	7 (7.93)	4 (6.29)	49
	+10%	3 (–8.77)	0	2 (–1.14)	3 (–1.13)	14	2 (0.91)	2 (1.83)	26
	+20%	0	0	0	0	5 (–2.60)	23	6 (3.27)	34
	+30%	2 (–4.50)	4 (–8.37)	5 (–16.21)	5 (–5.75)	3 (–1.81)	16 (–3.83)	76	111

and the average deviation in the expected R&D output. The results are depicted in Table 4:

We observe that the average expected R&D output is higher in group A compared to the one in group B (127.33 versus 122.85), and the average deviation of expected R&D output in group A is lower than the one in group B (7.84 versus 7.96). Based on this analysis, though it is not statistically significant, the robust model uses a policy of forming projects with high average and low deviation of expected R&D output in order to maximize expected R&D output when uncertainty is given.

## 6. Discussion

Determining R&D budget allocation is an important issue for many funding agencies. Though they have continually improved their budget allocation systems, there is still room for improvement. In parallel with their efforts, we devise a hybrid framework to allocate R&D budgets to maximize the expected R&D output, which is a common metric used by well-known research funding agencies [59]. Specifically, our study provides an innovative way of allocating R&D budgets by combining a state-of-the-art machine learning algorithm and an optimization technique. Considering that the purpose of our study is to design a systematic framework for evaluating a project's success probability and allocating budgets accordingly, the proposed framework fits our purpose. We also believe that the proposed framework acts as a decision support system for project managers to systematically select and allocate budgets to maximize the expected R&D output, as well as to minimize the risk of accepting projects. Moreover, it is worth noting that ours is formulated in a generic manner. Hence, it can be applied to various R&D budget allocation problems by additionally considering detailed constraints. This means that our framework can be used as a baseline of such extensions.

We find that we verify the importance of applying a systematic budget allocation scheme. Specifically, experimental results confirm that the funded projects under the regular funding scheme present lower performance than those selected using our approach. A reasonable interpretation of this result is that the studied funding agency performs poorly at picking “winners.” For example, the majority of selecting and funding protocols are designed in a qualitative manner (e.g., a panel review system) as many funding agencies use. Though such a system has been shown to be effective for evaluating projects, problems have been raised. Typically, prior studies have noted that such a system is weak at addressing bias [60,61]. For instance, an evaluation by an expert would be subjective, and may be politically

biased. The qualification of an expert who should be assigned to an appropriate knowledge to evaluate projects can also be a critical issue.

Because we provide a data-driven decision support framework for funding agencies to support objective and unbiased evaluations when determining budget allocation decisions, we believe that it plays a complementary role to the current protocols used by funding agencies. Moreover, considering the trend from emerging countries such as China that the number of public R&D programs has increased, the legacy system, a peer review system, can demand excessive time from experts, which may lead to the degradation of system's quality. Therefore, our method provides another benefit to minimize the burden of the evaluation process, as well as improving efficiency of the project selection and budget allocation process.

Furthermore, our study shows the added importance of considering uncertainty in allocating R&D budgets. That is, we verify the so-called *price of robustness*, which is the deterioration of the optimal value because owing the robust optimization model excludes nonrobust solutions to protect against any deviations in uncertain parameters. However, as seen in Fig. 4, simulation results with 100 randomly generated instances revealed that robust solutions outperform nominal solutions. Specifically, by choosing projects with higher average performance and lower deviation of the expected R&D output, robust solutions provide better performance than that of the nominal solutions. In other words, at least in our experimental settings, our model provides a budget allocation plan that shows solid performance against uncertainty in the expected R&D output.

Based on our research findings, our model can be used to enhance the overall quality of R&D budget allocation protocols. That is, it can be used as a useful tool for funding agencies to achieve objectives such as maximizing the expected R&D output. In particular, it can play a crucial role in selecting projects that are expected to achieve high-quality R&D outputs under the uncertainty of R&D outputs.

Our work sheds light on the problem of systematically planning public R&D budgets. However, there are a few limitations to our study. First, though we devise a hybrid framework combining an R&D output estimation model with an optimization model, each model can be improved. Specifically, using AutoML provides a convenient way of using state-of-the-art machine learning algorithms; however, it lacks explainability. In other words, if AutoML selects deep learning as the winning algorithm, we cannot know the importance of each independent variable because the calculation of deep learning is technically unknown. Second, our budget allocation model provides a single-period plan by considering a performance measure only (total sum of publications and patents). Though it can be a good starting point, many practical cases require more realistic constraints such as the planning period. Also, as noted earlier, though our problem instances are solved within a reasonable time, computational issues should be addressed to expand the applicability of our framework.

Despite the practical limitations, this study has substantial implications because it proposes an approach based on mathematical models for the optimal distribution of government R&D budgets. The approach is verified using an experiment and actual data.

**Table 4**  
Average of the expected R&D outputs and deviations in two different groups.

Metric	Group A	Group B
Average of the expected R&D output	127.33	122.85
Average deviations of the expected R&D output	7.84	7.96

## 7. Conclusion

This study addresses the issue of determining an optimal budget allocation for a national R&D program. For this, we propose a decision support framework to allocate the optimal amount of a budget for each project (within a specific R&D program) to maximize the sum of expected R&D outputs. In particular, we first use machine learning algorithms to estimate the expected R&D output for each project, and apply a robust optimization method to deal with the uncertainty in the expected output. Even though there are studies on R&D budget allocation and on project selection, a hybrid combination of prediction and robust optimization is lacking.

Moreover, we verify the actual applicability of the proposed framework by examining the BK21 Plus program, which is one of the large programs organized by the Korean Ministry of Education. In particular, we increased the validity of our work by using R&D-related information that can be obtained at the national level. As a result, we practically confirm the importance of budget allocation efficiency and demonstrate the price of robustness. Specifically, we reveal that when budget allocation efficiency is achieved, 13.6% better R&D output can be obtained, even using the same budget. This has significance as a case that shows the possibility of using a data-drive decision support framework in R&D budget compilation and allocation. Furthermore, we show the effectiveness of the robust model when uncertainty in the expected R&D output is given. Specifically, by assuming the realized deviation in the expected R&D output, the budget allocation plans obtained from the robust model outperform the one from the nominal model. We believe this highlights the practical value of our study. To sum up, our study establishes the foundation to contribute to the de-

velopment of relevant study areas by proposing a new decision support framework as well as by verifying the effectiveness of the proposed framework using nationwide collected data of government R&D program.

A few areas need to be considered in future studies. One such area could be technical advancement in the use of machine learning. For instance, explainability should be investigated further. Using state-of-the-art machine learning algorithms requires that we interpret the results reasonably. However, many machine learning algorithms are regarded as a black box, which means decision-makers may hesitate to use machine learning algorithms in practice. In addition, future efforts toward obtaining a generalizable estimation model should be conducted. Many extensions based on our current budget allocation model could be studied. For instance, because many R&D programs are typically conducted over many years, considering such a property when allocating R&D budgets is needed. Specifically, considering that a large number of R&D projects run for several years, some funding agencies may want to adjust the budget of certain projects based on the interim evaluation results. In that case, a research question of designing an optimal funding policy to maximize the long-term expected R&D output is required by considering the temporal aspect of R&D projects.

## Acknowledgments

This work is partially supported by the research project conducted in Science and Technology Policy Institute (grant number P0181800). The author is especially grateful to Dr. Woo, Chungwon, and Kim, Tae Kyung for their suggestions and insights, which have greatly improved this manuscript.

## Appendix A. Appendix

### A.1. Computational results for larger problem sizes

We have arbitrarily increased the problem size by virtually creating problem instances and have examined whether computational burden is observed. The computational results are presented in Table 5. As seen in Table 5, we cannot find significant computational load when problem sizes are enlarged up to 300% of the original volume.

However, it should be noted that as our problem is formulated as binary programming, the problem size directly affects complexity. This may cause a computational issue, particularly for larger problem cases.

**Table 5**  
CPU time in 11 scenarios.

Problem instances (Incremental rate of $B$ )	CPU Time (sec)		
	Original	Doubled (= 200% of its original volume)	Tripled
–20%	34.74	25.15	27.23
–15%	27.13	20.27	23.47
–10%	31.42	24.42	30.30
–5%	30.91	32.31	35.09
0%	39.20	31.00	33.28
+5%	40.42	28.26	29.18
+10%	34.21	30.19	34.12
+15%	27.28	32.13	34.52
+20%	35.62	19.13	21.81
+25%	28.44	24.58	26.99
+30%	30.01	31.23	33.84

### A.2. Sensitivity analysis

In Section 3, we simply added the number of publications and patents as a single unit to predict future R&D output. Though simple, some R&D funding agencies may want to assign different weights to the number of publications ( $w_{pub}$ ) and patents ( $w_{pat}$ ). Accordingly, we extended our original model ( $P^0$ ) by considering two different measures. To this end, we first calculate weighted values of two R&D outputs ( $e_{ij}'$ ). Then, by using this value, we predict future R&D output. Then, by using the extended model shown below ( $P^{ext}$ ), we find optimal R&D budget allocation plans.

$$e_{ij}' = w_{pub} e_{ij}^{pub} + w_{pat} e_{ij}^{pat} \quad (16)$$

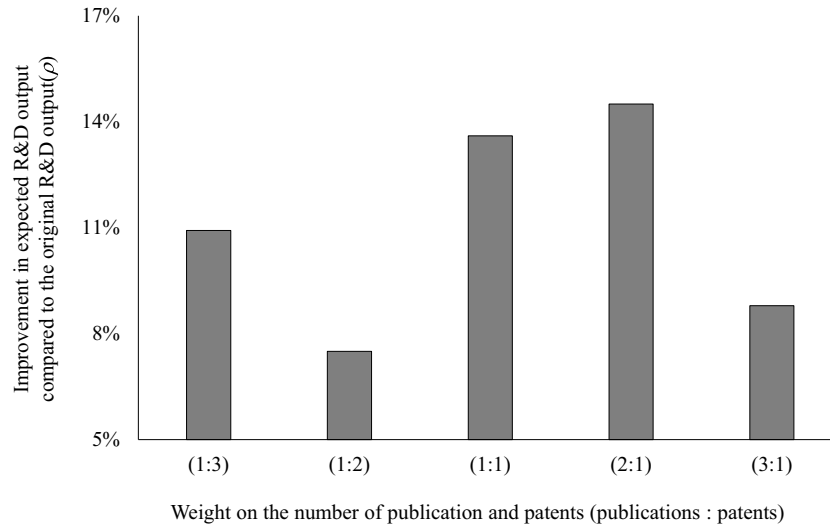


Fig. 5. Improvement in R&D output when different weights on the number of publications and patents are applied.

In  $(P^{ext})$ , we replace  $e_{ij}$  with  $e'_{ij}$ . All the others are same as in  $(P^0)$ . The test was conducted using a baseline scenario in which the ratio of  $w_{pub}$  to  $w_{pat}$  varied from (1:3) to (3:1). In total, we examine five different instances. Note that (1:1) is equivalent to the ratio presented in Section 3.

$$(P^{ext}) \max \sum_{i \in I} \sum_{j \in J} (e'_{ij}) y_{ij} \quad (17)$$

subject to Eqs. (7)–(9).

For each instance, we evaluate the value of  $\rho$ . The results are shown as follows (Fig. 5): Overall, the average of estimated performance improvement is computed by 11.06% in five problem instances. Though the magnitude of the performance improvement varied among the weight values, this analysis clearly verifies the effectiveness of using the proposed framework, regardless of considering different weights for the performance metrics.

## References

- [1] H. Park, J. Lee, B. Kim, Project selection in NIH: a natural experiment from ARRA, *Research Policy* 44 (2015) 1145–1159.
- [2] M. Brown, A. Forsythe, Robust tests for the equality of variances, *Journal of the American Statistical Association* 69 (1974) 364–367.
- [3] J. Fedderke, M. Goldschmidt, Does massive funding support of researchers work?: evaluating the impact of the South African research chair funding initiative, *Research Policy* 44 (2015) 467–482.
- [4] C. Grimpe, Extramural research grants and scientists' funding strategies: beggars cannot be choosers? *Research Policy* 41 (2012) 1448–1460.
- [5] H. Jang, C. Woo, T. Kim, A Study for Designing Optimal R&D Portfolios, Report from Science and Technology Policy Institute, Sejong, Republic of Korea, 2018.
- [6] A. Salo, J. Keisler, A. Morton, Portfolio Decision Analysis, *International Series in Operations Research and Management Science*, Springer.
- [7] F. Fabozzi, P.N. Kolm, D.A. Pachamanova, S.M. Focardi, Robust Portfolio Optimization and Management, Wiley, Hoboken, 2007.
- [8] P. Xidonas, G. Mavrotas, T. Krintas, J. Psarras, C. Zopounidis, Multicriteria Portfolio Management, Springer, New York, 2012.
- [9] R. Balachandra, J.H. Friar, Factors for success in R&D projects and new product innovation: a contextual framework, *IEEE Transactions on Engineering Management* 44 (1997) 276–287.
- [10] J. Cho, J. Lee, Development of a new technology product evaluation model for assessing commercialization opportunities using Delphi method and fuzzy AHP approach, *Expert Systems with Applications* 40 (2013) 5314–5330.
- [11] L.A. Heslop, E. McGregor, M. Griffith, Development of a technology readiness assessment measure: the cloverleaf model of technology transfer, *The Journal of Technology Transfer* 26 (2001) 369–384.
- [12] C. Liu, A study for allocating resources to research and development programs by integrated fuzzy DEA and fuzzy AHP, *Scientific Research and Essays* 6 (2011) 3973–3978.
- [13] H. Eilat, B. Golany, A. Shtub, R&D project evaluation: an integrated DEA and balanced scorecard approach, *Omega* 36 (2008) 895–912.
- [14] M. Talias, Optimal decision indices for R&D project evaluation in the pharmaceutical industry: Pearson Index versus Gittins Index, *European Journal of Operational Research* 177 (2007) 1105–1112.
- [15] C. Galbraith, A. DeNoble, S. Ehrlich, D. Kline, Can experts really assess future technology success? A neural network and Bayesian analysis of early stage technology proposals, *The Journal of High Technology Management Research* 17 (2007) 125–137.
- [16] B. Tan, E. Anderson Jr., J. Dyer, G. Parker, Evaluating system dynamics models of risky projects using decision trees: alternative energy projects as an illustrative example, *System Dynamics Review* 26 (2010) 1–17.
- [17] I. Seo, D. Lee, Analysis of effects on national productivity by R&D investment policy: focus on simulation results via system dynamics, *Journal of Governmental Studies* 16 (2010) 91–122.
- [18] T. Walworth, M. Yearworth, J. Davis, P. Davies, Early Estimation of Project Performance: The Application of a System Dynamics Rework Model, *IEEE International Systems Conference (SysCon)*, (2013).
- [19] X.H. Jin, G. Zhang, Modelling optimal risk allocation in PPP projects using artificial neural networks, *International Journal of Project Management* 29 (2011) 591–603.
- [20] Y. Wang, C. Yu, H. Chan, Predicting construction cost and schedule success using artificial neural networks ensemble and support vector machines classification models, *International Journal of Project Management* 30 (2012) 470–478.
- [21] Y. Wang, G. Gibson Jr., A study of preproject planning and project success using ANNs and regression models, *Automation in Construction* 19 (2010) 341–346.
- [22] Y. Chen, Y. Zhang, J. Liu, P. Mo, Interrelationships among critical success factors of construction projects based on the structural equation model, *Journal of Management in Engineering* 28 (2012).
- [23] D. Dvir, A. Ben-David, A. Sadeh, A. Shenhar, Critical managerial factors affecting defense projects success: a comparison between neural network and regression analysis, *Engineering Applications of Artificial Intelligence* 19 (2006) 535–543.
- [24] H. Song, H. Park, J. Sim, J. Seo, Ripple Effect Analysis of Government R&D Investment by System Approach, Report from Korea Institute of S&T Evaluation and Planning (2014).
- [25] J.H. Lorie, L.J. Savage, Three problems in rationing capital, *Journal of Business* 18 (1955) 229–239.
- [26] H. Markowitz, Portfolio selection, *Journal of Finance* 7 (1952) 77–91.
- [27] R.H. Bernhard, Mathematical programming models for capital budgeting—a survey, generalization, and critique, *Journal of Financial and Quantitative Analysis* 4 (1969) 111–158.
- [28] K. Heidenberger, C. Stummer, Research and development project selection and resource allocation: a review of quantitative modelling approaches, *International Journal of Management Reviews* 1 (1999) 197–224.
- [29] F. Ghasemzadeh, N.P. Archer, Project portfolio selection through decision support, *Decision Support Systems* 29 (2000) 73–88.
- [30] M. Badri, D. Davis, D. Davis, A comprehensive 0–1 goal programming model for project selection, *International Journal of Project Management* 19 (2001) 243–252.
- [31] S. Coldrick, P. Longhurst, P. Ivey, J. Hannis, An R&D options selection model for investment decisions, *Technovation* 25 (2005) 185–193.
- [32] S.A. Gabriel, S. Kumar, J. Ordóñez, A. Nasserian, A multiobjective optimization model for project selection with probabilistic considerations, *Socio-Economic Planning Sciences* 40 (2006) 297–313.

- [33] K. Doerner, W.J. Gutjahr, R.F. Hartl, C. Strauss, C. Stummer, Pareto ant colony optimization: a meta-heuristic approach to multi-objective portfolio selection, *Annals of Operations Research* 131 (2004) 79–99.
- [34] F. Abdelaziz, R. El Fayedh, A. Rao, A discrete stochastic goal program for portfolio selection: the case of United Arab Emirates equity market, *Information Systems and Operational Research* 47 (2009) 5–13.
- [35] A.L. Medaglia, D. Hueth, J. Mendieta, J. Sefair, A multiobjective model for the selection and timing of public enterprise projects, *Socio-Economic Planning Sciences* 42 (2008) 31–45.
- [36] A. Tolga, Fuzzy multicriteria R&D project selection with a real options valuation model, *Journal of Intelligent Fuzzy Systems* 19 (2008) 359–371.
- [37] A. Carazo, T. Gómez, J. Molina, A.G. Hernández-Díaz, M. Guerrero, R. Caballero, Solving a comprehensive model for multiobjective project portfolio selection, *Computers and Operations Research* 37 (2010) 630–639.
- [38] P. Guo, J. Liang, Y.M. Zhu, J.F. Hu, R&D project portfolio selection model analysis within project interdependencies context, *IEEE International Conference on Industrial Engineering and Engineering Management*, 2008.
- [39] W. Gutjahr, S. Katzensteiner, P. Reiter, C. Stummer, M. Denk, Competence-driven project portfolio selection, scheduling and staff assignment, *Central European Journal of Operations Research* 16 (2008) 281–306.
- [40] I. Litvinchev, F. López, A. Alvarez, E. Fernández, Large-scale public R&D portfolio selection by maximizing a biobjective impact measure, *IEEE Transactions on Systems, Man, and Cybernetics - Part A: Systems and Humans* 40 (2010) 572–582.
- [41] U. Jung, D.W. Seo, An ANP approach for R&D project evaluation based on interdependencies between research objectives and evaluation criteria, *Decision Support Systems* 49 (2010) 335–342.
- [42] M.B. Aryanezhad, H. Malekly, M. Karimi-Nasab, A fuzzy random multi-objective approach for portfolio selection, *Journal of Industrial Engineering International* 7 (2011) 12–21.
- [43] R. Bhattacharyya, S. Kar, D.D. Majumder, Fuzzy mean-variance-skewness portfolio selection models by interval analysis, *Computers & Mathematics with Applications* 61 (2011) 126–137.
- [44] L.M. Luo, Optimal diversification for R&D project portfolios, *Scientometrics* 91 (2012) 219–229.
- [45] S. Casault, A.J. Groen, J.D. Linton, Selection of a portfolio of R&D projects, Chapter 4, *Handbook on the Theory and Practice of Program Evaluation*, 2013, pp. 89–111.
- [46] F. Costantino, G.D. Gravio, F. Nonino, Project selection in project portfolio management: an artificial neural network model based on critical success factors, *International Journal of Project Management* 33 (2015) 1744–1754.
- [47] N.M. Arratia-Martínez, F. López, S.E. Schaeffer, L. Cruz-Reyes, Static R&D project portfolio selection in public organizations, *Decision Support Systems* 84 (2016) 53–63.
- [48] A. Soyster, Convex programming with set-inclusive constraints and applications to inexact linear programming, *Operations Research* 21 (1973) 1154–1157.
- [49] A. Ben-Tal, A. Nemirovski, Robust solutions of linear programming problems contaminated with uncertain data, *Mathematical Programming* 88 (2000) 411–424.
- [50] D. Bertsimas, M. Sim, The price of robustness, *Operations Research* 52 (2004) 35–53.
- [51] M.T. Melo, S. Nickel, F. Saldanha-da-Gama, Facility location and supply chain management—a review, *European Journal of Operational Research* 196 (2009) 401–412.
- [52] N. Gülpınar, D. Pachamanova, A robust optimization approach to asset-liability management under time-varying investment opportunities, *Journal of Banking & Finance* 37 (2013) 2031–2041.
- [53] K. Govindan, M. Fattahi, E. Keyvanshokoh, Supply chain network design under uncertainty: a comprehensive review and future research directions, *European Journal of Operational Research* 263 (2017) 108–141.
- [54] F. Hassanzadeh, H. Nemati, M. Sun, Robust optimization for interactive multi-objective programming with imprecise information applied to R&D project portfolio selection, *European Journal of Operational Research* 238 (2014) 41–53.
- [55] P. Mild, J. Liesiö, A. Salo, Selecting infrastructure maintenance projects with robust portfolio modeling, *Decision Support Systems* 77 (2015) 21–30.
- [56] S. Bekiros, J.A. Hernandez, S. Hammoudeh, D.K. Nguyen, Multivariate dependence risk and portfolio optimization: an application to mining stock portfolios, *Research Policy* 46 (2015) 1–11.
- [57] F. Liu, W.D. Zhu, Y.W. Chen, D.L. Xu, J.B. Yang, Evaluation, ranking and selection of R&D projects by multiple experts: an evidential reasoning rule based approach, *Scientometrics* 111 (2017) 1501–1519.
- [58] B. Sun, Y. Liu, G. Yang, A robust pharmaceutical R&D project portfolio optimization problem under cost and resource uncertainty, *Journal of Uncertain Systems* 11 (2017) 205–220.
- [59] E. Vilkkumaa, A. Salo, J. Liesiö, A. Siddiqui, Fostering breakthrough technologies—how do optimal funding decisions depend on evaluation accuracy? *Technological Forecasting and Social Change* 96 (2015) 173–190.
- [60] K. Boudreau, E. Guinan, K. Lakhani, C. Riedl, Looking across and looking beyond the knowledge frontier: intellectual distance, novelty, and resource allocation in science, *Management Science* 62 (2016) 2765–2783.
- [61] J. Wang, R. Veugelers, P. Stephan, Bias against novelty in science: a cautionary tale for users of bibliometric indicators, *Research Policy* 46 (2017) 1416–1436.

**Hoon Jang** is an associate research fellow at the Science and Technology Policy Institute (STEPI), Korea. He received a Ph.D. in Industrial & Systems Engineering from KAIST in 2014. His research interests lie in the area of data-driven system analytics using optimization methods. For the past few years, he has been working on designing a better R&D budgeting systems in Korea.