# Design of Experiments Algorithm for Comprehensive Exploration and Rapid Optimization in Chemical Space

**Kazuhiro Takeda[a]\*, Masaru Kondo[b], Muthu Karuppasamy[c,d], Mohamed S. H. Salem[c,e], and Shinobu Takizawa[c]**

[a] Shizuoka University, Department of Applied Chemistry and Biochemical Engineering, Hamamatsu, Shizuoka, Japan
[b] University of Shizuoka, Yada, Shizuoka, Japan
[c] SANKEN, Osaka University, Ibaraki, Osaka 5670047, Japan
[d] Graduate School of Pharmaceutical Sciences, Osaka University, 1-6 Yamada-oka, Suita, Osaka 565-0871, Japan
[e] Suez Canal University, Pharmaceutical Organic Chemistry Department, Faculty of Pharmacy, Ismailia 41522, Egypt
\* Corresponding Author: takeda.kazuhiro@shizuoka.ac.jp.

## ABSTRACT

Bayesian optimization is known to be able to search for the optimal conditions based on a small number of experiments. However, these experiments are insufficient to understand the experimental condition space. In contrast, we report the development of an algorithm that combines a low-confounding definitive screening design with Bayesian optimization, allowing for rapid optimization and ensuring sufficient experiments to understand the experimental condition space with a low confounding.

**Keywords**: Optimization, Algorithms, Bayesian optimization, Definitive screening design.

## INTRODUCTION

It has been established that Bayesian optimization (BO) [1,2] is capable of searching for the optimal conditions based on a limited number of experiments. Accordingly, prior to implementing the BO, preliminary experiments are conducted in accordance with the principles of the design of experiments. Many methods for the designs of experiment have been proposed, including the factorial design[3], the Plackett Burman design [4], the Box Behnken design [5], the central composite design[6], the Latin hypercube design [7], and the definitive screening design (DSD) [8]. The results of experiments designed by the DSD may contribute to an understanding of the chemical space with minimal confounding factors and a limited number of steps. Subsequently, preliminary experiments are conducted based on the design of experiment, utilizing the design of experiment methods, and a search for the optimal conditions is then undertaken by the BO. This enables the experimental space to be optimized following a comprehensive investigation and to gain insight into the surrounding conditions of the optimal parameters. However, as the design of experiment is conducted independent of the BO, it contributes no more than the preliminary experiments to the optimization
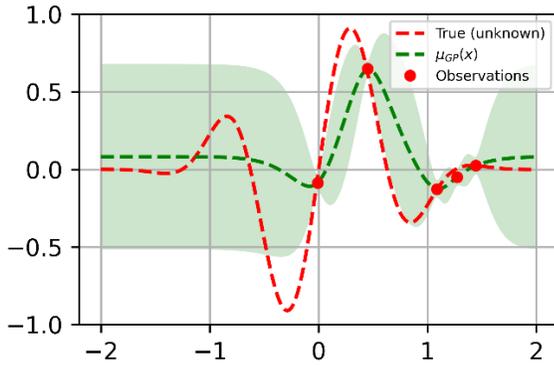
process. As a result of comparing the Box Behnken design and response surface methodology with adaptive Bayesian optimization, the BO should be combined with the design of experiments, since the BO could not investigate the experimental space [9]. Therefore, in this study, the design of experiment was evaluated by the acquisition function from within the DSD, and the design of experiment is accordingly modified. The results of running the modified design of experiment are suitable for optimization by the BO, thus allowing the subsequent BO to rapidly search for the optimum conditions.

This paper presents the development of an algorithm that combines a small confounding of the DSD with the BO, thus allowing for rapid optimization and ensuring sufficient experiments to understand the experimental condition space with a small confounding.
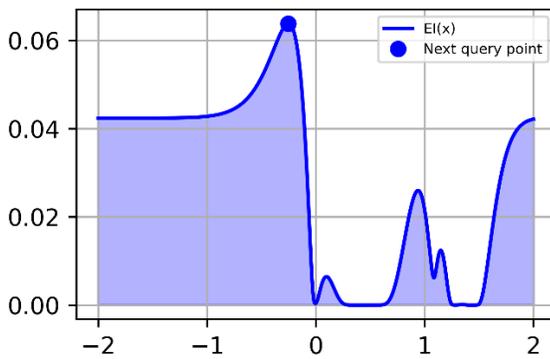
## BAYESIAN OPTIMIZATION

The BO is a method that is known to converge towards the optimum value based on a limited number of experiments. A schematic representation of the process by which the subsequent trial point is formulated is provided in Figure 1. A regression is performed on the observed values with the base estimator frequently being a

Gaussian process regression. An acquisition function is calculated based on the regression model, and the point with the maximum acquisition function value is the next trial point. The acquisition functions most commonly employed are LCB (lower confidence bound), EI (negative expected improvement), PI (negative probability of improvement) and gp_hedge, which probabilistically combines them.



(a) Gaussian process regression from observations.



(b) Next query point by acquisition function.

**Figure 1.** Schematics of Bayesian optimization.

## DEFINITIVE SCREENING DESIGN

The DSD [1,2] is one of the famous design methods of experiments. The DSD is characterised by the following attributes:

- Low confounding between factors.
- A relatively limited number of experiments.

$$N = 2n + 1 + 4f \qquad (1)$$

where $N$ is the number of experiments, $n$ is the number of factors, and $f$ is the number of fake factors.

- Three levels for each factor.

The fake factor allows the error to be estimated independently of model selection. The design of experiment with the DSD utilizing a 3 variables model without a fake factor is presented in Table 1. There are $7 (= 2 \cdot 3 + 1 + 4 \cdot 0)$ entries for 3 variables and no fake factor.

**Table 1:** Designed experiments by the DSD with 3 variables and no fake factor.

| Entry | A | B | C |
|-------|-----|-----|-----|
| 1 | 0 | 1 | 1 |
| 2 | 0 | -1 | -1 |
| 3 | 1 | 0 | -1 |
| 4 | -1 | 0 | 1 |
| 5 | 1 | -1 | 0 |
| 6 | -1 | 1 | 0 |
| 7 | 0 | 0 | 0 |

## INTEGRATION OF DEFINITIVE SCREENING DESIGN AND BAYESIAN OPTIMIZATION

The algorithm proposed in this study is presented in the following section. The initial stage of the study comprises preliminary experiments, which are conducted in accordance with the design of experiment formulated by the DSD (Step 1). However, the design of experiment is modified with the objective of incorporating the experimental data into the subsequent BO process. Specifically, in order to calculate the acquisition function for each factor, a preliminary experiment is conducted in accordance with the planned methodology until each variable value is no longer unique (Step 2). Subsequently, the acquisition function is calculated for the remaining experiments, and the experiment with the highest acquisition function value is selected for experimentation (Steps 3, 4, and 5). This process is repeated for Steps 3, 4, and 5 until no further designed experiments remain (Step 6). Once all the designed experiments have been conducted, the optimal conditions can be identified through Bayesian optimization (Steps 7, 8, and 9). By prioritizing the condition with the largest acquisition function value in Steps 3, 4, and 5, the experimental values that are most suitable for the Bayesian optimization performed after Step 7 can be acquired. Furthermore, the planned experimental conditions based on the DSD satisfy the DSD characteristics of a small confounding and a small number of experiments to understand the experimental space.

Example) The proposed algorithm is applied to the designed experiment in Table 1. Step 1 is to formulate Table 1. Step 2 is to execute the experiment up to Entry 3 as designed. Step 3 is to calculate the acquisition function values from Entry 4 onwards. In step 3, the acquisition function values for entry 4 and subsequent entries are calculated. In order to preserve the sign, the

acquisition function values are calculated in the range 0 to 1 for the condition value 1, in the range 0 to -1 for the condition value -1, and in the range 0 for the condition value 0. In Step 6, Steps 3 to 5 are executed until no further entries remain. In Step 7, the condition with the maximum acquisition function over the entire range is identified based on the experimental results up to Step 6, and experiments are conducted using Step 8. In Step 9, Steps 7 and 8 are repeated until convergence is confirmed.

**Proposed Algorithm**

Step 1: Design experiment by the DSD.

Step 2: Experiment until all elements are not unique under the DSD definition, in which + and – are the maximum and minimum between the chemical space, respectively.

Step 3: Determine all the remaining experimental conditions for the maximum acquisition function value within the same sign region.

Step 4: Select the experimental condition with the maximum acquisition function value between those of the remaining designed experiments.

Step 5: Perform an experiment with the determined conditions.

Step 6: Go to the next step if all the experiments are executed, else return to Step 3.

Step 7: Search the optimum condition for the maximum acquisition function value.

Step 8: Perform an experiment with the condition.

Step 9: Finish if the convergence conditions are satisfied, else go to Step 7.

## NUMERICAL SIMULATION

In order to ascertain the efficacy of the algorithm proposed in this study, a numerical comparison was conducted. In the present study, the design of experiment was formulated by the DSD and compared with nine different designs of experiment using PyDOE2 [10]. These included a factorial design, a Plackett Burman design, a Box Behnken design, and the DSD. The DSD had three types of designs of experiment (0, 2, and 4) depending on the number of fake factors. Furthermore, a comparison was conducted between the proposed algorithm and Bayesian optimization with randomly generated initial values and no design of experiment. In order to examine the differences between the design of experiment and Bayesian optimization in terms of how they are combined, the following cases were considered: In the case of

'Sequential BO', the BO was performed subsequent to the design of the experiment. In the case of 'In BO', the experimental conditions were set to the condition with the maximum acquisition function value based on the results of Step 3. In the case of 'In Remain BO', the experimental conditions were set to the condition with the maximum acquisition function value based on the acquisition function value and the order of the design with the latter being switched based on the acquisition function values, as determined by Step 4.
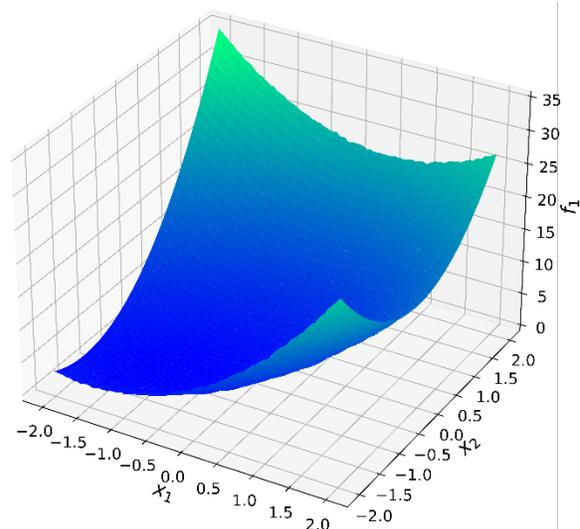


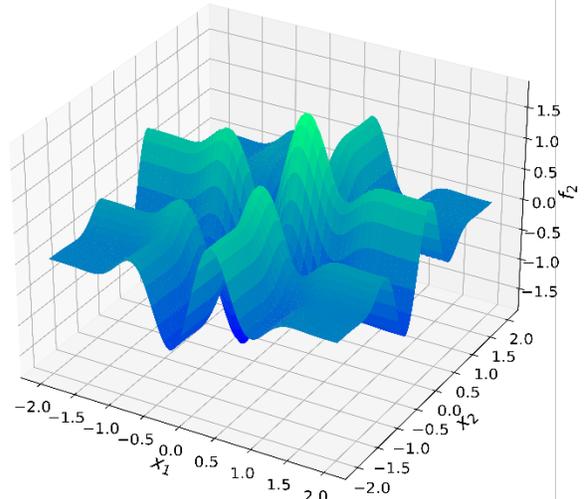**Figure 2.** Quadratic objective function with one peak.



**Figure 3.** Tanh objective function with many peaks.

The objective functions were two distinct types: unimodal and multimodal. The former is illustrated in Eq. (2) (Figure 2), while the latter is shown in Eq. (3) (Figure 3). In Figure 3, $x_3 = 0$ for illustration.

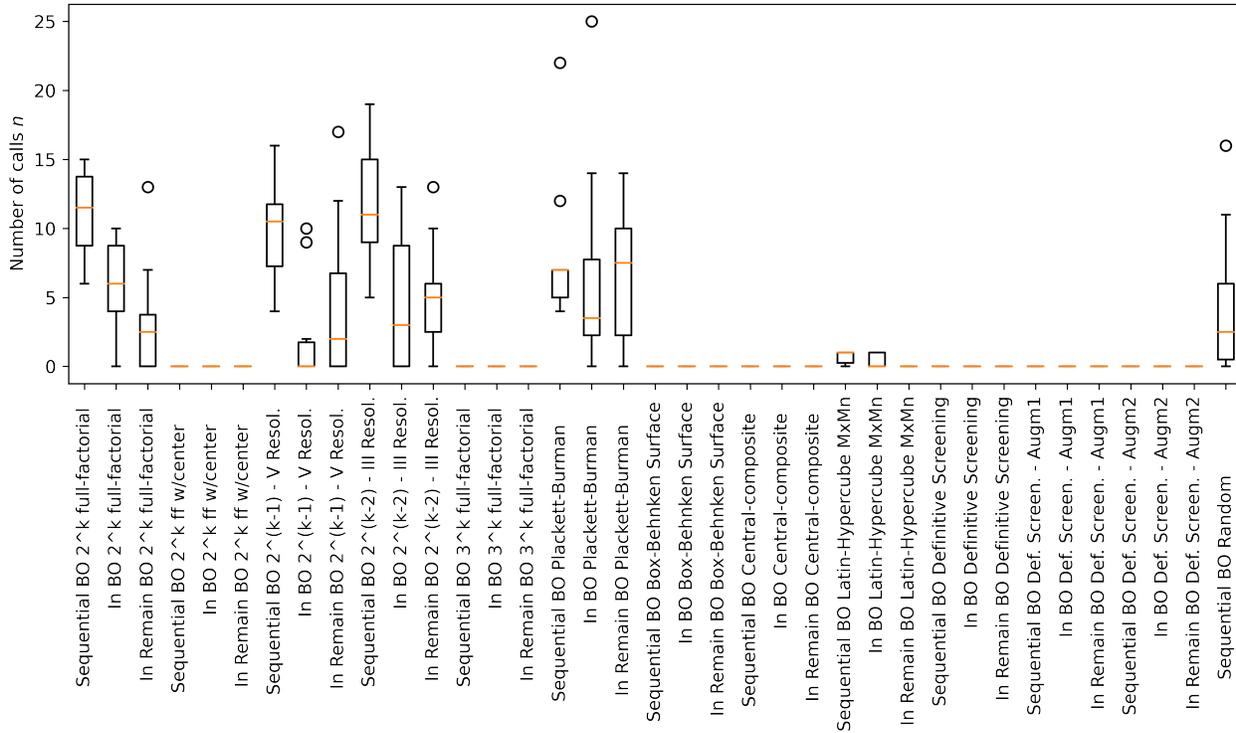$$f_1(x_1, x_2) = 3 + 2x_1 + 4x_2 - 2x_1^2 - 2x_1 x_2 + 3x_2^2 \quad (2)$$

**Figure 4:** Box plots for number of iterations for quadratic function using each algorithm. Orange line represents 50%tile, box represents from 25%tile to 75%tile, whiskers represent minimum and maximum values, and circle represents outliers. For 'Sequential BO', the BO was executed after DoE. For 'In BO', DoE was affected by acquisition function. For 'In Remain BO', DoE was affected and rearranged by acquisition function. '2^k full-factorial' is 2 level general full factorial, '2^k ff w/center' is 2 level general full factorial with zero matrix, '2^(k-1) - V Resol.' Is V resolution 2 level fractional factorial, '2^(k-2) - III Resol.' Is III resolution 2 level fractional factorial, '3^k full-factorial' is 3 level general full factorial, 'Plackett-Burman' is Plackett Burman, 'Box-Behnken Surface' is Box Behnken, 'Central-composite' is central composite, 'Latin-Hypercube MxMn' is Latin hypercube, 'Definitive Screening' is the DSD, 'Def. Screen. - Augm1' is the DSD with 2 fake factors, and 'Def. Screen. - Augm2' is the DSD with 4 fake factors. 'Random' is random initial conditions.

$$f_2(x_1, x_2, x_3) = \sin(5x_1) \cdot (1 - \tanh(x_1^2))$$
$$+ \sin(5x_2) \cdot (1 - \tanh(x_2^2))$$
$$+ \sin(5x_3) \cdot (1 - \tanh(x_3^2)) \quad (3)$$

The functions were constrained as follows.

$$-2.0 \le x_1 \le 3.0, -3.0 \le x_2 \le 4.0, -2.0 \le x_3 \le 3.0 \quad (4)$$

A total of 10 trials were conducted for each case in the simulation.

## RESULTS AND DISCUSSION

### Case 1: unimodal objective function

Figure 4 illustrates the distribution of experiments conducted to optimize the unimodal objective function. The number of experiments designed by each method was quite different, the Plackett Burman required only 4 experiments but the 3 level general full factorial required 27 experiments. As the total number of experiments including the designed experiments was not suitable for evaluating the performance of Bayesian optimization after designed experiments, the number of calls after each designed experiment was counted. As the number of convergences of the BO varied between experiments, ten experiments were conducted until a neighbourhood of the optimum value of 3.0 was reached at 5.0. As the optimum value of the unimodal objective function is readily discernible, in numerous instances the experimental conditions proximate to the optimum value were attained during the preliminary experimental phase. The results demonstrated that convergence was achieved in a significantly reduced number of experiments compared to those conducted using the 'Sequential BO Random' method, in which the initial values were randomly generated.

In the case of the 2-level general full factorial, designated as a '2^k full-factorial', and the III resolution 2-
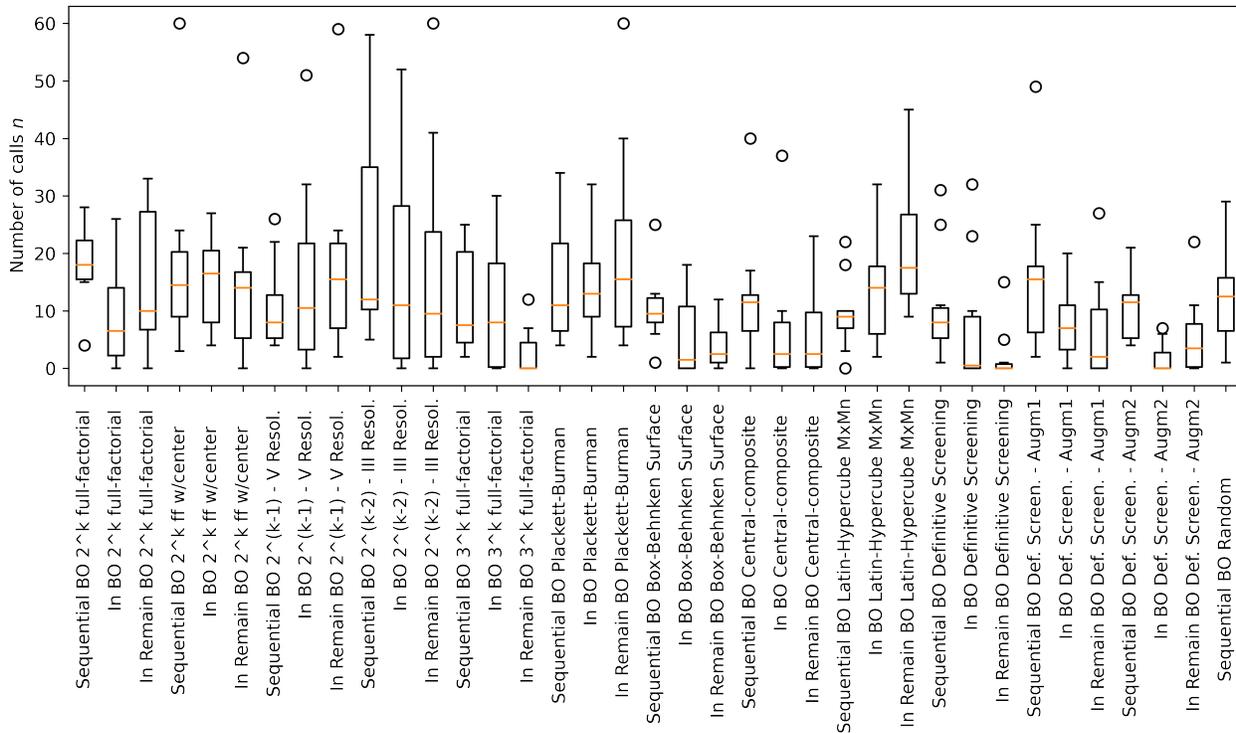
**Figure 5:** Box plots for number of iterations for quadratic function using each algorithm. Orange line represents 50%tile, box represents from 25%tile to 75%tile, whiskers represent minimum and maximum values, and circle represents outliers. For 'Sequential BO', the BO was executed after DoE. For 'In BO', DoE was affected by acquisition function. For 'In Remain BO', DoE was affected and rearranged by acquisition function. '2^k full-factorial' is 2 level general full factorial, '2^k ff w/center' is 2 level general full factorial with zero matrix, '2^(k-1) - V Resol.' Is V resolution 2 level fractional factorial, '2^(k-2) - III Resol.' Is III resolution 2 level fractional factorial, '3^k full-factorial' is 3 level general full factorial, 'Plackett-Burman' is Plackett Burman, 'Box-Behnken Surface' is Box Behnken, 'Central-composite' is central composite, 'Latin-Hypercube MxMn' is Latin hypercube, 'Definitive Screening' is the DSD, 'Def. Screen. - Augm1' is the DSD with 2 fake factors, and 'Def. Screen. - Augm2' is the DSD with 4 fake factors. 'Random' is random initial conditions.

level fractional factorial, designated as '2^(k-2) - III Resol.', it is evident that the case of 'In BO' converged. This occurred less frequently than in the cases of the 'Sequential BO' and 'In Remain BO'. The fact that the case of 'In BO' converged less often than the case of 'Sequential BO' demonstrates the value of optimizing the experimental conditions through the acquisition function during the preliminary experiment based on the design of experiment. This demonstrated the value of optimizing the experimental conditions with the acquisition function during the preliminary experiment based on the design of experiment. The case of the 'In Remain BO' converged fewer times than the case of 'In BO', indicating the benefit of utilizing the acquisition function to alter the design order during the preliminary experiment based on the design of experiment.

## Case 2: multimodal objective function

Figure 5 illustrates the distribution of experiments conducted for the purpose of optimizing the multimodal objective function. Similarly, for the unimodal function, to evaluate the performance of Bayesian optimization after designed experiments, the number of calls after each designed experiment was counted. Ten experiments were conducted for the multimodal objective function with the number of experiments continuing until a neighborhood of the optimum value of -2.8 was reached at -1.0. In comparison to the unimodal objective function, the search for the optimal value of the multimodal objective function was more challenging and necessitated the execution of several tens of experiments.

In a considerable number of designs of experiment, the case of 'In BO' converged less often than the case of 'Sequential BO', and the case of 'Sequential BO' converged less often than the case of 'In Remain BO'

In comparison to the 'Sequential BO Random' method, which employed randomly generated initial values, i.e., the Box-Behnken 'Box-Behnken Surface' and the

central Composite 'Central-composite' methods, the DSD 'Definitive Screening', and the 'Def. The following abbreviations were used: 'Screen', 'Augm1' and 'Def'. Furthermore, the 'Def. Augm1' and 'Def. Augm2' methods were also considered. The 'Screen. - Augm2' method demonstrated a high degree of convergence in a relatively small number of cycles. Notably, in the case of employing the proposed algorithm within the DSD framework, it exhibited a markedly superior convergence rate compared to the 'Sequential BO Random' method.

In both cases, the proposed algorithm, in which the DSD designed experiments were influenced and rearranged by the acquisition function, could converge faster than the 'Sequential BO Random' method. These findings highlight the value of optimizing the experimental conditions through the acquisition function, and of modifying the design of experiment order through the acquisition function during the preliminary experiment, independent of the objective function. Since faster convergence for BO means that the exploration of the experimental condition space was more useful. Therefore, in the proposed algorithm, the BO was able to perform on sufficient experiments designed by the DSD to understand the experimental condition space.

## CONCLUSIONS

We put forth a novel algorithm that integrates the DSD and the BO to address confounding factors and optimize the experimental conditions. Optimizing the magnitude and orders of the designed experiments of the DSD based on the BO contributed to effectively accelerate the convergence rate. Our proposed algorithm exhibits the potential for a faster convergence than existing algorithms in numerical experiments. Furthermore, we have demonstrated the value of rearranging the order of experiments in achieving more efficient outcomes. Furthermore, in the proposed algorithm, the BO was able to perform on sufficient experiments designed by the DSD to understand the experimental condition space.

In future work, we will investigate wet experiments to verify the efficacy of our proposed algorithm.

## ACKNOWLEDGEMENTS

## REFERENCES

1. Snoek J., Larochelle H., Adams R. P., Practical Bayesian Optimization of Machine Learning Algorithms, *arXiv*:1206.2944:1-9 (2012) https://doi.org/10.48550/arXiv.1206.2944
2. Kondo M., Wathsala H. D. P., Ishikawa K., Yamashita D., Miyazaki T., Ohno Y., Sasai H., Washio T., Takizawa S., Bayesian optimization-assisted screening to identify improved reaction conditions for Spiro-Dithiolane synthesis, *Molecules*, 28, 13, 5180 (2023) https://doi.org/10.3390/molecules28135180
3. Fisher, R. A., The arrangement of field experiments. *Journal of the Ministry of Agriculture*. 33:503-515 (1926) https://doi.org/10.23637/rothamsted.8v61q
4. Plackett, R. L., Burman, J. P., The Design of Optimum Multifactorial Experiments. *Biometrika*, 33(4):305–325 (1946) https://doi.org/10.2307/2332195
5. Box, G. E. P., Behnken, D. W., Some New Three Level Designs for the Study of Quantitative Variables. *Technometrics*, 2(4):455–475 (1960) https://doi.org/10.1080/00401706.1960.10489912
6. Box,G.E.P., Wilson,K.B., On the Experimental Attainment of Optimum Conditions, *Journal of Royal Statistical Society Series B*, 13:1-45 (1951) https://doi.org/10.1111/j.2517-6161.1951.tb00067.x
7. Mckay, M. D., Beckman, R. J., Conover, W. J., A Comparison of Three Methods for Selecting Values of Input Variables in the Analysis of Output from a Computer Code. *Technometrics*, 42(1):55–61 (2000) https://doi.org/10.1080/00401706.2000.10485979
8. Jones, B., Nachtsheim, C. J., A Class of Three-Level Designs for Definitive Screening in the Presence of Second-Order Effects. *Journal of Quality Technology*, 43(1):1–15 (2011) https://doi.org/10.1080/00224065.2011.11917841
9. Rummukainen H., Horhammer H., Kuusela P., Kilpi J., Sirvio J., Makela M., Traditional or adaptive design of experiments? A pilot-scale comparison on wood delignification. *Heliyon*, 10:e24484 (2024) https://doi.org/10.1016/j.heliyon.2024.e24484
10. Sjogren R., Svensson D., https://github.com/clicumu/pyDOE2