

Modeling and Active Learning for Experiments with Quantitative-Sequence Factors

Qian Xiao¹, Yaping Wang², Abhyuday Mandal¹, and Xinwei Deng³



¹ University of Georgia, Athens, GA ² East China Normal University, Shanghai, China ³ Virginia Tech, Blacksburg, VA

Quantitative-Sequence (QS) factors:

• In modern scientific areas, there are non-traditional experiments considering both the **quantities** and **sequences** for arranging components, named as quantitative-sequence (QS) factors.

• Cancer Treatment – in vitro study:

- Three anti-tumor drugs A , B and C were added in a sequence with different doses.
- The percentage of tumor inhibition was measured six hours after administering the last drug.

Run	Drug A		Drug B		Drug C		Response
	dosage	order	dosage	order	dosage	order	
1	3.75 μM	1	95 nM	2	0.16 μM	3	39.91
2	2.80 μM	1	70 nM	2	0.16 μM	3	30.00
3	3.75 μM	3	95 nM	1	0.16 μM	2	34.68

• Experiments with QS Inputs

– Characteristics of Quantitative-Sequence (QS) factor:

- * Different drug dosage affects the response.
- * Different sequence order affects the response.
- * The QS factor is not purely continuous, not purely categorical and not purely ordinal.

– Objectives:

- * To study the relationship between the response and the QS factor.
- * To optimize the dosage and the order in the sequence for each drug.

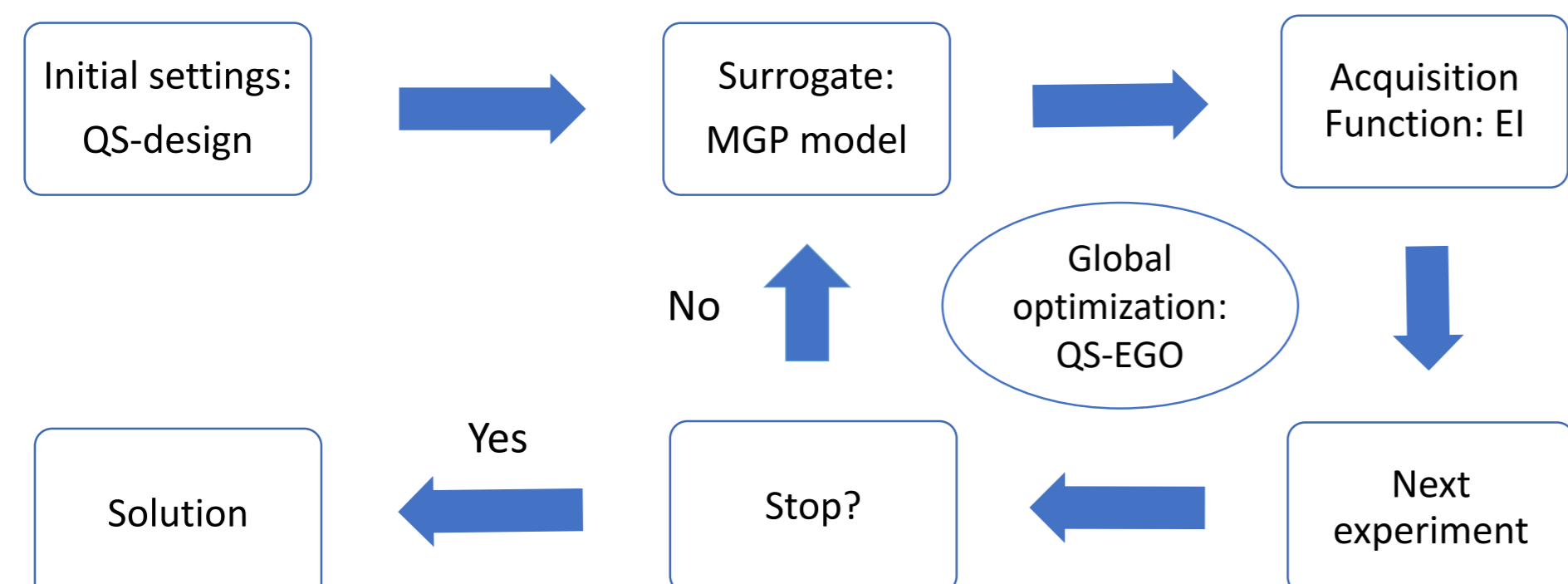
– Challenges:

- * Numerous possibilities: for k components with s levels, $s^k \times k!$ possible runs.
- * A good design for QS input is not trivial.
- * A good statistical model for QS input is needed.

QS-learning

We propose an active learning approach (QS-learning) which includes

1. **MaGP**: a novel mapping-based additive Gaussian process model for prediction and uncertainty quantification,
2. **QS-EGO**: a sequential scheme using efficient global optimization algorithms,
3. **QS-design**: a new class of optimal experimental designs for collecting initial data points.



MaGP

- Consider the i^{th} input as $\mathbf{w}_i = (\mathbf{x}_i^T, \mathbf{o}_i^T)^T$ with \mathbf{x}_i takes quantitative values and \mathbf{o}_i is a vector containing the orders of the components in the arrangement sequence.
- For an experiment with n runs and k components, we model the output at $\mathbf{w} = (\mathbf{x}^T, \mathbf{o}^T)^T$ as

$$Y(\mathbf{w}) = \mu + \sum_{h=1}^k G_h(\mathbf{w}) + \epsilon,$$

where G_1, \dots, G_k are independent zero-mean GP with stationary covariance functions and $\epsilon \sim N(0, \tau^2)$ is a random error ($\tau^2 > 0$ for physical experiments and $\tau^2 = 0$ for computer experiments).

- G_h corresponds to the impact of h^{th} component with its covariance function ϕ_h as

$$\phi_h(\mathbf{w}_i, \mathbf{w}_j) = \sigma_h^2 \exp\left\{-\theta_h(x_{i,h} - x_{j,h})^2\right\} \exp\left\{-\sum_{l=1}^t (\delta_{i,h}^{(l)} - \delta_{j,h}^{(l)})^2\right\}$$

where σ_h^2 is the variance parameter and θ_h is the correlation parameter for the h^{th} component.

- We consider mapping the order o_h to a vector $(\delta_h^{(1)}, \dots, \delta_h^{(t)})$. The t -dimensional mapping ($t = 1, \dots, k-1$) for the order of any component is defined as

$$\begin{bmatrix} c_1, \dots, c_k \\ 1 \\ 2 \\ \vdots \\ k \end{bmatrix} \rightarrow \begin{bmatrix} \delta_1^{(1)} & \delta_1^{(2)} & \dots & \delta_1^{(t)} \\ \delta_2^{(1)} & \delta_2^{(2)} & \dots & \delta_2^{(t)} \\ \vdots & \vdots & \vdots & \vdots \\ \delta_k^{(1)} & \delta_k^{(2)} & \dots & \delta_k^{(t)} \end{bmatrix}_{k \times t},$$

where we set $\delta_l^{(j)} = 0$ for all $j \geq l$ to avoid over-parametrization.

QS-EGO

1. **Step 1**: Construct an optimal initial design for QS factors with n_0 runs $\mathbf{w}_1, \dots, \mathbf{w}_{n_0}$, evaluate their responses $Y(\mathbf{w}_1), \dots, Y(\mathbf{w}_{n_0})$, and fit the MaGP model based on these observations.

2. **Step 2**: Let the next design point \mathbf{w}_{n+1} maximize the expected improvement $E[I(\mathbf{w}_{n+1})]$ and observe $Y(\mathbf{w}_{n+1})$. (We proposed an algorithm.)

Here for a target input $\mathbf{w}_* = (\mathbf{x}_*, \mathbf{o}_*)$:

$$E[I(\mathbf{w}_*)] = (y_{\min}^{(n)} - \hat{Y}(\mathbf{w}_*))\Phi\left(\frac{y_{\min}^{(n)} - \hat{Y}(\mathbf{w}_*)}{s(\mathbf{w}_*)}\right) + s(\mathbf{w}_*)\varphi\left(\frac{y_{\min}^{(n)} - \hat{Y}(\mathbf{w}_*)}{s(\mathbf{w}_*)}\right)$$

3. **Step 3**: Re-fit the MaGP model based on observations $(\mathbf{w}_1, Y(\mathbf{w}_1)), \dots, (\mathbf{w}_{n+1}, Y(\mathbf{w}_{n+1}))$.

4. **Step 4**: Repeat Step 2 and 3 until the stopping criterion is met or the maximum number of sequential runs is reached.

QS-design

1. We propose a new class of optimal designs for QS factors, named as **QS-design**, which achieves **space-filling** and **pair-balanced** properties.

2. We propose a general approach to construct QS-design with any run and factor sizes; and provide a deterministic algebraic construction for certain design sizes.

3. Denote the design for QS factors as $D = (X, O)$ where X is the quantitative part and O is the sequence part, both using components as columns.

4. Our key idea is to first construct a good design O , and then construct a good design X in combination with O to obtain the QS-design $D = (X, O)$.

A Real Combinatorial Drug Experiment on Lymphoma

Lymphoma is cancer that causes lymphocytes grow out of control.

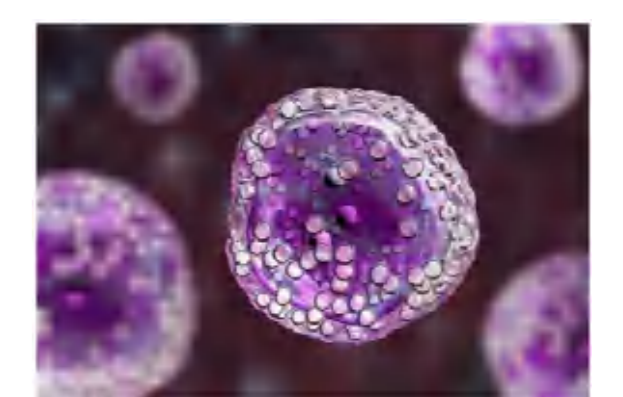
- Real Data: a **24-run** three-drug (A : paclitaxel, B : doxorubicin, C : mitoxantrone) experiment for Lymphoma cancer treatment. (Wang et al. 2020)

– All six possible sequences of the three drugs were enumerated. For each sequence, two dose-levels for A (Level 0: 2.8 μM ; Level 1: 3.75 μM) and B (Level 0: 70 nM; Level 1: 95 nM), and a fixed dose-level for C (0.16 μM) were considered.

– We run the proposed QS-learning to see if we can use fewer runs to identify the optimal treatment in this experiment. We construct an **8-run QS-design** to collect the initial data and the proposed QS-learning selects **7 sequential runs** until the stopping rule is satisfied.

– The true maximum response **47.18** has been found, along with the third and fourth largest responses **44.38** and **44.33**.

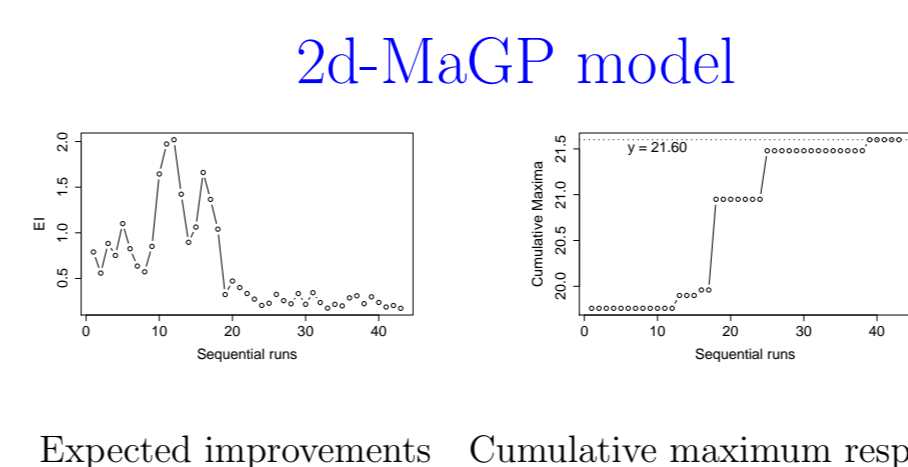
X_A	O_A	X_B	O_B	O_C	Y
0	3	0	2	1	35.04
1	2	0	1	3	22.26
0	1	1	3	2	43.93
1	2	1	3	1	20.88
0	1	0	2	3	30.00
1	1	0	3	2	38.18
0	2	1	1	3	26.02
1	3	1	1	2	34.68



Simulation Studies

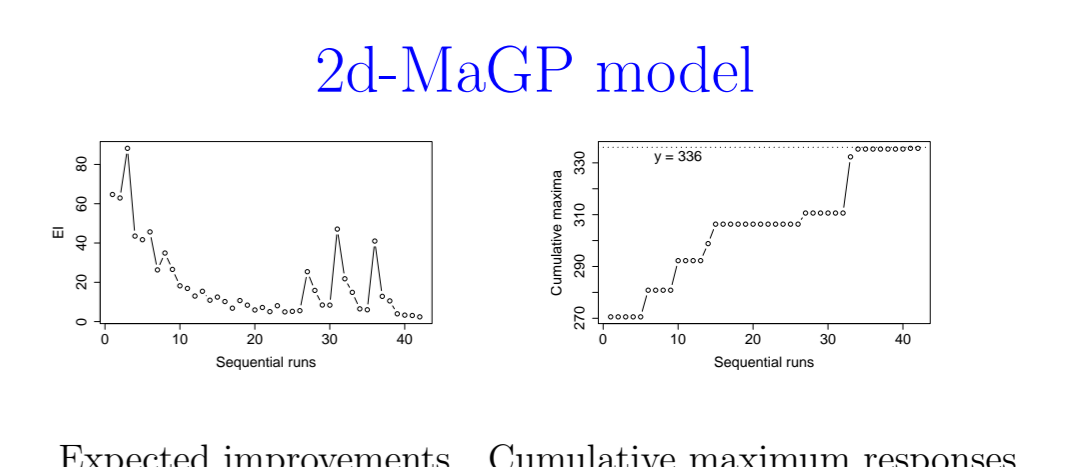
- Single Machine Scheduling Problem

$$R(\mathbf{x}) - C(\boldsymbol{\alpha}, \mathbf{x}) = w_0 \sum_{i=1}^k x_i - \sum_{h=1}^k w_h T^2(\alpha_h)$$



- Traveling Salesman Problem

$$F(\mathbf{x}, \boldsymbol{\alpha}) = ka + e \sum_{i=1}^k x_i - bC(\mathbf{x}, \alpha_k) - f \sum_{j=1}^k T(\mathbf{x}, \alpha_j).$$



Conclusions

- In this work, we propose an active learning approach to identify optimal solutions for experiments with quantitative-sequence (QS) factors.
- Analyzing such experiments is challenging due to their semi-discrete and possibly very large solution spaces as well as complex input-output relationships.
- From our empirical results, the proposed QS-learning can provide desirable solutions within a few number of sequential runs.