

Problem Set 01

Published on 11.04.2019

To be collected on 16.04.2019

Problem 1: (20 points)

This problem comes from a famous method called High Dimensional Model Representation (HDMR) or ANalysis Of VAriance (ANOVA). A simple definition of the HDMR or ANOVA data decomposition model is as follows:

Let \mathbb{I} denote the unit interval $[0, 1]$, \mathbb{I}^n the n -dimensional unit hypercube and the random variable $\mathbf{x} \in \mathbb{I}^n$. Let $f : \mathbb{I}^n \rightarrow \mathbb{R}$ represent a general dynamic system, which depends on the n -D random variable $\mathbf{x} = [x_1, \dots, x_n]^T$. Suppose that the system function $f(\mathbf{x})$ is integrable in a reasonable certain of degree, e.g. square integrable or third-order integrable, it could be represented as a sum of 2^n sub-components defined in all possible sub-spaces:

$$f(\mathbf{x}) = f_0 + \sum_{d=1}^n \sum_{i_1 < \dots < i_d} f_{(i_1, \dots, i_d)}(x_{i_1}, \dots, x_{i_d}), \quad (1)$$

where the interior sum is over all sets of d integers i_1, \dots, i_d that satisfy $1 \leq i_1 < \dots < i_d \leq n$. Relation (1) also can be understood as

$$f(\mathbf{x}) = f_0 + \sum_i f_{(i)}(x_i) + \sum_{i < j} f_{(i,j)}(x_i, x_j) + \dots + f_{(1,2,\dots,n)}(x_1, x_2, \dots, x_n). \quad (2)$$

For example, taking $n = 3$, we have $\mathbf{x} = [x_1, x_2, x_3]^T \in \mathbb{I}^3$, and

$$f(\mathbf{x}) = f_0 + f_{(1)}(x_1) + f_{(2)}(x_2) + f_{(3)}(x_3) + f_{(1,2)}(x_1, x_2) + f_{(1,3)}(x_1, x_3) + f_{(2,3)}(x_2, x_3) + f_{(1,2,3)}(x_1, x_2, x_3). \quad (3)$$

The first summand f_0 is the global mean of $f(\mathbf{x})$, e.g.

$$f_0 := \int_{\mathbb{I}^n} f(\mathbf{x}) d\mathbf{x} := \int_{\mathbb{I}} \dots \int_{\mathbb{I}} f(\mathbf{x}) dx_1 \dots dx_n, \quad (4)$$

while other summands can be obtained recursively

$$\begin{aligned} \int_{\mathbb{I}^{n-1}} f(\mathbf{x}) \prod_{k \neq i} dx_k &= f_0 + f_{(i)}(x_i), \quad i = 1, 2, \dots, n \\ \int_{\mathbb{I}^{n-2}} f(\mathbf{x}) \prod_{k \neq i, j} dx_k &= f_0 + f_{(i)}(x_i) + f_{(j)}(x_j) + f_{(i,j)}(x_i, x_j), \quad i \neq j = 1, 2, \dots, n \\ &\dots = \dots \\ f(\mathbf{x}) &= f_0 + \sum_{d=1}^{n-1} \sum_{i_1 < \dots < i_d} f_{(i_1, \dots, i_d)}(x_{i_1}, \dots, x_{i_d}) + f_{(1,2,\dots,n)}(x_1, x_2, \dots, x_n). \end{aligned} \quad (5)$$

If we reconsider the example shown in (3), we can obtain each summand from the given $f(\mathbf{x})$, e.g.

$$\begin{aligned} f_{(2)}(x_2) &:= \int_{\mathbb{I}} \int_{\mathbb{I}} f(\mathbf{x}) dx_1 dx_3 - f_0 \\ f_{(1,3)}(x_1, x_3) &:= \int_{\mathbb{I}} f(\mathbf{x}) dx_2 - f_0 - f_{(1)}(x_1) - f_{(3)}(x_3). \end{aligned} \quad (6)$$

Now, if we define the full sub-index set $\mathbb{S} := \{1, 2, \dots, n\}$, its subset $\mathbb{U} \subseteq \mathbb{S}$, and the size of the set $|\mathbb{U}|$, e.g. $|\mathbb{S}| = n$, the whole HDMR or ANOVA can be compactly represented as

$$\begin{aligned} f(\mathbf{x}) &= \sum_{\mathbb{U} \subseteq \mathbb{S}} f_{\mathbb{U}}(\mathbf{x}_{\mathbb{U}}), \\ f_{\mathbb{U}}(\mathbf{x}_{\mathbb{U}}) &= \int_{\mathbb{I}^{n-|\mathbb{U}|}} f(\mathbf{x}) d\mathbf{x}_{\mathbb{S} \setminus \mathbb{U}} - \sum_{\mathbb{V} \subset \mathbb{U}} f_{\mathbb{V}}(\mathbf{x}_{\mathbb{V}}). \end{aligned} \quad (7)$$

Now, with the basic definitions in the probability and statistics, please prove following nice properties:

(1) Exclude f_\emptyset , all other summands $f_U(\mathbf{x}_U)$ has the zero mean (7 points)

$$\int_{\mathbb{I}} f_U(\mathbf{x}_U) d\mathbf{x}_k = 0, \quad U \subseteq (\mathbb{S} \setminus \emptyset), \quad k \in U; \quad (8)$$

(2) Include f_\emptyset , every two distinct summands $f_U(\mathbf{x}_U)$ and $f_V(\mathbf{x}_V)$ are orthogonal (7 points)

$$\int_{\mathbb{I}^n} f_U(\mathbf{x}_U) f_V(\mathbf{x}_V) d\mathbf{x} = 0, \quad U, V \subseteq \mathbb{S}, \quad U \neq V; \quad (9)$$

(3) Define the variance of each summand and the total variance of $f(\mathbf{x})$ as

$$D_U := \int_{\mathbb{I}^{|U|}} f_U^2(\mathbf{x}_U) d\mathbf{x}_U, \quad U \subseteq (\mathbb{S} \setminus \emptyset) \quad (10)$$

$$D := \int_{\mathbb{I}^n} (f(\mathbf{x}) - f_\emptyset)^2 d\mathbf{x}$$

the following equality exists (6 points)

$$D := \sum_{U \subseteq (\mathbb{S} \setminus \emptyset)} D_U. \quad (11)$$

(Additional Bonus: 10 points)

If we define the third-order moment, the skewness, and the fourth-order moment, the kurtosis, as follows

$$\gamma := \int_{\mathbb{I}^n} (f(\mathbf{x}) - f_\emptyset)^3 d\mathbf{x}, \quad (12)$$

$$\kappa := \int_{\mathbb{I}^n} (f(\mathbf{x}) - f_\emptyset)^4 d\mathbf{x},$$

Can you derive the corresponding equality between γ and γ_U , and the one between κ and κ_U ?

Remarks:

The reason that I changed my original mind and design such big problem as the only problem in set 01 is:

- (a) this is a very famous data decomposition method based on the concepts in statistics, which is unfortunately not introduced in many statistical data analysis textbooks;
- (b) if we divide the both sides of (11) by D , and define that

$$1 := \sum_{U \subseteq (\mathbb{S} \setminus \emptyset)} \frac{D_U}{D} =: \sum_{U \subseteq (\mathbb{S} \setminus \emptyset)} S_U, \quad (13)$$

the S_U could be understood as the sensitivity index of those random variables \mathbf{x}_U to the dynamic system $f(\mathbf{x})$. The larger S_U , the more sensitive of the system to \mathbf{x}_U . In other words, those \mathbf{x}_U will be more important to the system. In the future, if you are facing to an unfamiliar big system or black box, and would like to understand it or control it, hope that you shall take some benefits from here.

Two useful references:

- [1] I.M. Sobol, Theorems and examples on high dimensional model representation, Reliability Engineering and System Safety, vol. 79, pp. 187-193, 2003.
- [2] M. Holtz, Sparse grid quadrature in high dimensions with applications in finance and insurance, Chapter 2, Springer-Verlag Berlin Heidelberg, 2011.