

Measuring the interactions among variables of functions over $[0, 1]^n$

Jean-Luc Marichal and Pierre Mathonet

University of Luxembourg

Measure of influence of variables

Problem

Let $f: [0, 1]^n \rightarrow \mathbb{R}$

We want to measure the *influence* (*importance*) of x_k over

$$y = f(x_1, \dots, x_k, \dots, x_n)$$

Example: affine function

$$y = c_0 + c_1 x_1 + \dots + c_k x_k + \dots + c_n x_n$$

Measure of influence of x_k :

coefficient c_k

Measure of influence of variables

More complicated functions...

$$y = \prod_{i=1}^n x_i^{c_i}$$

$$y = \left(\sum_{i=1}^n c_i x_i^2 \right)^{1/2}$$

$$y = \max_{i=1..n} \min(c_i, x_i)$$

etc.

Measure of influence of variables

A reasonable answer:

→ Approximation of f by an affine function

$$y = a_0 + a_1 x_1 + \cdots + a_k x_k + \cdots + a_n x_n$$

Measure of influence of x_k over f :

coefficient a_k

Measure of interaction among variables

What about interactions among variables?

Example: multilinear function

$$f(x_1, x_2) = c_0 + c_1 x_1 + c_2 x_2 + c_{12} x_1 x_2$$

Measure of *interaction* between x_1 and x_2 within f :

coefficient c_{12}

$c_{12} = 0$: zero interaction

$c_{12} > 0$: positive interaction

$c_{12} < 0$: negative interaction

Measure of interaction among variables

Let $S \subseteq N = \{1, \dots, n\}$

Measure of interaction among variables $\{x_k : k \in S\}$

→ Approx. of f by a multilinear polynomial of degree $\leq s = |S|$

$$f_s(\mathbf{x}) = \sum_{\substack{T \subseteq N \\ |T| \leq s}} a_s(T) \prod_{i \in T} x_i$$

Measure of interaction among $\{x_k : k \in S\}$ inside f :

coefficient $a_s(S)$

(leading coefficients)

Approximation problem

Multilinear approximation

Denote by M_s the set of all multilinear polynomials $g: [0, 1]^n \rightarrow \mathbb{R}$ of degree $\leq s$

$$g(\mathbf{x}) = \sum_{\substack{T \subseteq N \\ |T| \leq s}} a_s(T) \prod_{i \in T} x_i$$

We define the *best s-th approximation* of a function $f \in L^2([0, 1]^n)$ as the multilinear polynomial $f_s \in M_s$ that minimizes

$$d(f, g)^2 = \int_{[0,1]^n} (f(\mathbf{x}) - g(\mathbf{x}))^2 d\mathbf{x}$$

among all $g \in M_s$

Interaction index

From the best s -th approximation of f , we define the following *interaction index* (*power index* if $s = 1$)

$$\mathcal{I}(f, S) = a_s(S) \quad (s = |S|)$$

= coefficient of $\prod_{i \in S} x_i$ in the best s -th approximation of f

In the discrete case, i.e., $f: \{0, 1\}^n \rightarrow \mathbb{R}$

- Non-weighted distance: Hammer and Holzman (1992)
- Weighted (multiplicative) distance: M. and Mathonet (2010)

Interaction index

Explicit expression

$$\mathcal{I}(f, S) = 12^s \int_{[0,1]^n} f(\mathbf{x}) \prod_{i \in S} \left(x_i - \frac{1}{2} \right) d\mathbf{x}$$

Properties:

- $f \mapsto \mathcal{I}(f, S)$ is linear and continuous
- The index \mathcal{I} is symmetric
- $f \in M_n \Rightarrow \mathcal{I}(f, S) = I_B(f|_{\{0,1\}^n}, S)$ (Banzhaf index)
- ...

Question: How can we see that this index actually measures an interaction?

Interaction index

Interpretation:

$$\begin{aligned} D_k f(\mathbf{x}) &= \text{rate of change w.r.t. } x_k \text{ of } f \text{ at } \mathbf{x} \\ &= \text{local contribution of } x_k \text{ over } f \text{ at } \mathbf{x} \end{aligned}$$

$$D_j D_k f(\mathbf{x}) = \text{local interaction between } x_j \text{ and } x_k \text{ within } f \text{ at } \mathbf{x}$$

Theorem

If f is sufficiently smooth, then

$$\mathcal{I}(f, S) = \int_{[0,1]^n} q_S(\mathbf{x}) D_S f(\mathbf{x}) \, d\mathbf{x}$$

where $q_S(\mathbf{x})$ is the p.d.f. of independent beta distributions $(2, 2)$

Interaction index

What if f is not differentiable?

S -derivative \rightarrow *discrete S -derivative* (S -difference quotient)

Theorem

We have

$$\mathcal{I}(f, S) = \int_{\mathbf{x} \in [0,1]^n} \int_{\mathbf{h}_S \in [\mathbf{0}, \mathbf{1} - \mathbf{x}_S]} p_S(\mathbf{h}) \frac{\Delta_{\mathbf{h}}^S f(\mathbf{x})}{\prod_{i \in S} h_i} d\mathbf{h}_S d\mathbf{x}$$

where $p_S(\mathbf{h})$ is a p.d.f. over the domain of integration

Interaction index

Examples:

$$f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n x_i$$

- $\mathcal{I}(f, \{k\}) = \frac{1}{n}$
- $\mathcal{I}(f, S) = 0 \quad s \geq 2$

$$f(\mathbf{x}) = \left(\prod_{i=1}^n x_i \right)^{1/n}$$

- $\mathcal{I}(f, \{k\}) = \left(\frac{n}{n+1} \right)^n \frac{6}{n+2}$
- $\mathcal{I}(f, S) = \left(\frac{n}{n+1} \right)^n \left(\frac{6}{n+2} \right)^s$

Interaction index

Further examples:

We have explicit expressions for $\mathcal{I}(f, S)$ when

$$f(\mathbf{x}) = \sum_{T \subseteq N} c(T) \prod_{i \in T} \varphi_i(x_i)$$

(pseudo-multilinear function)

$$f(\mathbf{x}) = \sum_{T \subseteq N} c(T) \min_{i \in T} x_i$$

(Choquet integral)

Further properties

A variable x_k is *inefficient* for f if f does not depend on x_k

$$f(x_k; \mathbf{x}_{N \setminus k}) - f(0_k; \mathbf{x}_{N \setminus k}) = 0$$

Property

If x_k is inefficient for f then

$$\mathcal{I}(f, S) = 0 \quad \forall S \ni k$$

Discrete case: k null player

Further properties

A variable x_k is *dummy* for f if

$$f(x_k; \mathbf{x}_{N \setminus k}) - f(0_k; \mathbf{x}_{N \setminus k}) = f(x_k; \mathbf{0}_{N \setminus k}) - f(0_k; \mathbf{0}_{N \setminus k})$$

(the marginal contribution of x_k is independent of $\mathbf{x}_{N \setminus k}$)

Discrete case: k dummy player

Further properties

A combination of variables $\{x_k : k \in S\}$ is *dummy* for f if

$$f(\mathbf{x}_S; \mathbf{x}_{N \setminus S}) - f(\mathbf{0}_S; \mathbf{x}_{N \setminus S}) = f(\mathbf{x}_S; \mathbf{0}_{N \setminus S}) - f(\mathbf{0}_S; \mathbf{0}_{N \setminus S})$$

Property

If $\{x_k : k \in S\}$ is *dummy* for f , then

$$\mathcal{I}(f, T) = 0 \quad \forall T \text{ such that } T \cap S \neq \emptyset \text{ and } T \setminus S \neq \emptyset$$

Discrete case: dummy coalition

Conclusion

Given a function $f \in L^2([0, 1]^n)$, we have defined an index to measure

- the influence of variables over f
- the interaction among variables within f

Interpretations

- leading coefficients in multilinear approximation (multilinear regression)
- mean S -derivative or S -difference quotient

Properties

- linearity, symmetry, S -monotonicity...
- natural extension of the Banzhaf interaction index (null and dummy players...)

Thank you for your attention!

arXiv : 0912.1547