

1. Shapley value theorem:

We assume that the payoff function below is definition of the function of the Shapley value theorem:

$$\Phi(i, N) = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|! (|N| - |S| - 1)!}{|N|!} (v(S \cup \{i\}) - v(S)) \quad (1.1)$$

let R represent a permutation of players set N

let P_i^R represent the set of players permuted before player i in permutation R . then (1.1) is equivalent to:

$$\Phi(i, N) = \sum_R \frac{1}{|N|!} (v(P_i^R \cup \{i\}) - v(P_i^R)) \quad (1.2)$$

$$1). \text{ Symmetry: } v(S \cup \{i\}) = v(S \cup \{j\}), \text{ if } \forall S, i, j \notin S \quad (1.3)$$

Based on (1.2):

$$\Phi(i, N) = \sum_R \frac{1}{|N|!} (v(P_i^R \cup \{i\}) - v(P_i^R)) \quad (1.4)$$

$$\Phi(j, N) = \sum_R \frac{1}{|N|!} (v(P_j^R \cup \{j\}) - v(P_j^R)) \quad (1.5)$$

Let the permutation whose i before j be R_1 .

Let the permutation whose j before i be R_2 .

Knowably, $R = R_1 \cup R_2$ with $R_1 \cap R_2 = \emptyset$ (1.6)

for $\forall r \in R_1, r = [\underbrace{\dots}_{r_0}, \underbrace{i, \dots}_{r_1}, \underbrace{j, \dots}_{r_2}]$, correspondingly we have

and only have one $r' \in R_2, r' = [\underbrace{\dots}_{r_0}, \underbrace{j, \dots}_{r_1}, \underbrace{i, \dots}_{r_2}]$ with

the part-permutation r_0 , r_1 and r_2 are the same.

Obviously, $P_i^r = P_j^{r'}$, based on (1.3), we have:

$$v(P_i^r \cup \{i\}) - v(P_i^r) = v(P_j^{r'} \cup \{j\}) - v(P_j^{r'}) \quad (\text{one-to-one mapping between } r \text{ \& } r')$$

$$\Rightarrow \sum_{R_1} \frac{1}{|N|!} (v(P_i^{R_1} \cup \{i\}) - v(P_i^{R_1})) = \sum_{R_2} \frac{1}{|N|!} (v(P_j^{R_2} \cup \{j\}) - v(P_j^{R_2})) \quad (1.7)$$

$$\text{Similarly, } \sum_{R_2} \frac{1}{|N|!} (v(P_i^{R_2} \cup \{i\}) - v(P_i^{R_2})) = \sum_{R_1} \frac{1}{|N|!} (v(P_j^{R_1} \cup \{j\}) - v(P_j^{R_1})) \quad (1.8)$$

Based on (1.6), (1.4) and (1.5) can be transformed into

$$\begin{cases} \Phi(i, N) = \sum_{R_1} \frac{1}{|N|!} (v(P_i^{R_1} \cup \{i\}) - v(P_i^{R_1})) + \sum_{R_2} \frac{1}{|N|!} (v(P_i^{R_2} \cup \{i\}) - v(P_i^{R_2})) \quad (1.9) \\ \Phi(j, N) = \sum_{R_1} \frac{1}{|N|!} (v(P_j^{R_1} \cup \{j\}) - v(P_j^{R_1})) + \sum_{R_2} \frac{1}{|N|!} (v(P_j^{R_2} \cup \{j\}) - v(P_j^{R_2})) \quad (1.10) \end{cases}$$

$$\begin{cases} \Phi(i, N) = \sum_{R_1} \frac{1}{|N|!} (v(P_i^{R_1} \cup \{i\}) - v(P_i^{R_1})) + \sum_{R_2} \frac{1}{|N|!} (v(P_i^{R_2} \cup \{i\}) - v(P_i^{R_2})) \quad (1.9) \\ \Phi(j, N) = \sum_{R_1} \frac{1}{|N|!} (v(P_j^{R_1} \cup \{j\}) - v(P_j^{R_1})) + \sum_{R_2} \frac{1}{|N|!} (v(P_j^{R_2} \cup \{j\}) - v(P_j^{R_2})) \quad (1.10) \end{cases}$$

Bringing (1.7) & (1.8) into (1.9) respectively yields (1.10)

$\Rightarrow \Phi(i, N) = \Phi(j, N)$, Symmetry satisfied.

2. Dummy: $v(S \cup \{i\}) = v(S) + v(\{i\})$, $\forall S, i \notin S$ (1.11)

Based on (1.2) and (1.11)

$$\begin{aligned} \Phi(i, N) &= \sum_R \frac{1}{|N|!} (v(P_i^R \cup \{i\}) - v(P_i^R)) \\ &= \sum_R \frac{1}{|N|!} (v(P_i^R) + v(\{i\}) - v(P_i^R)) \\ &= \frac{1}{|N|!} v(\{i\}) \cdot n_R \quad (1.12) \end{aligned}$$

n_R is a full permutation of N , so $n_R = |N|!$ (1.13)

Bringing (1.13) into (1.12), we have:

$\Phi(i, N) = v(\{i\})$, Dummy satisfied

3) Additivity: $v = v_1 + v_2$ (1.14)

Based on (1.2) and (1.14):

$$\begin{aligned}
\Phi(i, N) &= \sum_R \frac{1}{|N|!} (V(P_i^R \cup \{i\}) - V(P_i^R)) \\
&= \sum_R \frac{1}{|N|!} [V_1(P_i^R \cup \{i\}) + V_2(P_i^R \cup \{i\}) - (V_1(P_i^R) + V_2(P_i^R))] \\
&= \sum_R \frac{1}{|N|!} [(V_1(P_i^R \cup \{i\}) - V_1(P_i^R)) + (V_2(P_i^R \cup \{i\}) - V_2(P_i^R))] \\
&= \sum_R \frac{1}{|N|!} (V_1(P_i^R \cup \{i\}) - V_1(P_i^R)) + \sum_R \frac{1}{|N|!} (V_2(P_i^R \cup \{i\}) - V_2(P_i^R)) \\
&= \Phi_1(i, N) + \Phi_2(i, N), \text{ Additively satisfied}
\end{aligned}$$

Because the payoff function (1.1) satisfies
 1). Symmetry, 2). Dummy 3) Additivity, it is
 the unique definition of the Shapley value

2. Integrated Gradient Theorem:

Definition of IG: $\frac{\partial f(\vec{x}) + \alpha(\vec{x} - \vec{x}')}{\partial [x_i + \alpha(x_i - x_i')]} \rightarrow$

$$\Phi(i, \vec{x}) = (x_i - x_i') \int_{\alpha=0}^1 \frac{\partial}{\partial x_i} f(\vec{x}' + \alpha(\vec{x} - \vec{x}')) d\alpha \quad (2.1)$$

1) Completeness: let Path $\vec{r}(\alpha) = \vec{x}' + \alpha(\vec{x} - \vec{x}')$ $\alpha \in [0, 1]$ (2.2)

Base on (2.1) and (2.2), we have:

$$\begin{aligned}
f(\vec{x}) - f(\vec{x}') &= f(\vec{r}(1)) - f(\vec{r}(0)) = \int_{\vec{r}(0)}^{\vec{r}(1)} \frac{d(f(\vec{r}(\alpha)))}{d(\vec{r}(\alpha))} d(\vec{r}(\alpha)) \\
&= \int_{\alpha=0}^1 \frac{d(f(\vec{r}(\alpha)))}{d(\vec{r}(\alpha))} \frac{d(\vec{r}(\alpha))}{d\alpha} d\alpha = \int_{\alpha=0}^1 \langle \nabla_{\vec{r}} f(\vec{r}(\alpha)), \vec{r}'(\alpha) \rangle d\alpha
\end{aligned}$$

$$\begin{aligned}
&(\vec{a} = [a_1, \dots, a_n]^T, \vec{\beta} = [b_1, \dots, b_n]^T, \langle \vec{a}, \vec{\beta} \rangle = \vec{a}^T \vec{\beta} = \sum_{i=1}^n a_i b_i) \\
&= \sum_i \int_0^1 [\nabla_{\vec{r}} f(\vec{r}(\alpha))]_i \cdot [\vec{r}'(\alpha)]_i d\alpha = \sum_i \int_0^1 \frac{\partial f(\vec{r}(\alpha))}{\partial [\vec{r}(\alpha)]_i} \cdot [x_i - x_i'] d\alpha \\
&= \sum_i [x_i - x_i'] \int_0^1 \frac{\partial}{\partial x_i} f(\vec{x}' + \alpha(\vec{x} - \vec{x}')) d\alpha = \sum_i \Phi(i, \vec{x})
\end{aligned}$$

$\Rightarrow f(\vec{x}) - f(\vec{x}') = \sum_i \phi(i, \vec{x}), (2.3)$ completeness satisfied!

2) linearity: $f = a_1 f_1 + a_2 f_2, (2.4)$ based on (2.3):

$$f_1(\vec{x}) - f_1(\vec{x}') = \sum_i \phi_1(i, \vec{x}), (2.5) \quad f_2(\vec{x}) - f_2(\vec{x}') = \sum_i \phi_2(i, \vec{x}), (2.6)$$

based on (2.4), (2.5), (2.6), we have:

$$\begin{aligned} \sum_i \phi(i, \vec{x}) &= f(\vec{x}) - f(\vec{x}') = a_1 f_1(\vec{x}) + a_2 f_2(\vec{x}) - a_1 f_1(\vec{x}') - a_2 f_2(\vec{x}') \\ &= a_1 f_1(\vec{x}) - a_1 f_1(\vec{x}') + a_2 f_2(\vec{x}) - a_2 f_2(\vec{x}') = a_1 \sum_i \phi_1(i, \vec{x}) + a_2 \sum_i \phi_2(i, \vec{x}) \\ &= \sum_i [a_1 \phi_1(i, \vec{x}) + a_2 \phi_2(i, \vec{x})] \end{aligned}$$

$\Rightarrow \phi(i, \vec{x}) = a_1 \phi_1(i, \vec{x}) + a_2 \phi_2(i, \vec{x}), (2.7)$ linearity satisfied!

3) Symmetry: f is Symmetric $\Rightarrow f(x_1, \dots, x_j, \dots, x_k, \dots, x_n)$
 $= f(x_1, \dots, x_k, \dots, x_j, \dots, x_n)$, for $\forall j \neq k \in [1, n]$ (2.8)

let $\vec{x}_{j,k} = (x_1, \dots, x_j, \dots, x_k, \dots, x_n)$, $\vec{x}'_{j,k} = (x'_1, \dots, x'_j, \dots, x'_k, \dots, x'_n)$

$\Rightarrow \vec{x}_{k,j} = (x_1, \dots, x_k, \dots, x_j, \dots, x_n)$, $\vec{x}'_{k,j} = (x'_1, \dots, x'_k, \dots, x'_j, \dots, x'_n)$

It is easy to prove: $f(\vec{x}_{j,k} + \alpha(\vec{x}'_{j,k} - \vec{x}_{j,k})) = f(\vec{x}_{k,j} + \alpha(\vec{x}'_{k,j} - \vec{x}_{k,j}))$ (2.9)

based on (2.1) & (2.9) we have:

$$\begin{aligned} 1^\circ \text{ for } \forall i \notin \{j, k\}, \phi(i, \vec{x}_{k,j}) &= (x_i - x'_i) \int_{\alpha=0}^1 \partial_{x_i} f(\vec{x}_{k,j} + \alpha(\vec{x}'_{k,j} - \vec{x}_{k,j})) d\alpha \\ &= (x_i - x'_i) \int_{\alpha=0}^1 \partial_{x_i} f(\vec{x}_{j,k} + \alpha(\vec{x}'_{j,k} - \vec{x}_{j,k})) d\alpha = \phi(i, \vec{x}_{j,k}) \end{aligned}$$

$$2^\circ \text{ for } i \in \{j, k\}, \phi(j, \vec{x}_{k,j}) = \phi(k, \vec{x}_{j,k}), \quad \phi(k, \vec{x}_{k,j}) = \phi(j, \vec{x}_{j,k})$$

\Rightarrow for \forall individual $\phi(i, \vec{x})$, symmetry satisfied

4) Sensitivity:

if $f(x)$ does not depend on i , then in (2.1)

$\partial_{x_i} f(x' + \alpha(x - x')) = 0 \Rightarrow \phi(i, x) = 0$, Sensitivity satisfied!

5). Implementation Invariance:

for different implementation of f , the math formula of f will keep its mapping relationship from each input and its corresponding output, so base on (2.1).

$\Phi(i, \vec{x})$ will just be related to the input-output mapping relationship of the numerical value of function. Implementation Invariance satisfied.

Because the payoff function (2.1) satisfies

- 1) Completeness 2) Linearity 3) Symmetry 4) Sensitivity
- 5) Implementation Invariance, it's the unique definition of payoff function of Integrated Gradient Theorem