

システムの変化に追従可能でかつ理解し易いドキュメントシステムのモデル化

1. モチベーション

- システムの運用をするためにはシステムの理解が必要
- ところが、システムは継続的に開発されるものであり変化するものである。また、その変化のスピードは早くなっている(要出典)
- また、システムは巨大で複雑なものになっており1人で全てを理解するのは現実的に難しくなっている
- 上記のような状況の中、継続的にシステムを運用していくためのシステム理解の効率的な手法を得たい
 - ゼロからシステムを理解する場合
 - 変化したシステムについて理解の追従をさせる場合

2. ターゲット

- 「既にユーザに機能を提供しており運用が必要なシステム」にターゲットを絞る。新規開発などの設計段階のシステムは含まない。
- システムを理解する手法としてドキュメントに注目する

3. システムとそれを説明するドキュメントのモデル化

前提

- エンジニアがシステムを理解しようとするとき、ドキュメントの内容が正しくシステムを説明している場合、システムを直接見て理解するより、ドキュメントを通じて理解するほうが理解し易いとする
 - つまり、本稿は「いかに正しいドキュメントを得るか」を主題とする

システム

- あるシステムにおいて、理解しなければならない要素を本稿では「システム要素」と呼ぶ。システム要素を s_1, s_2, \dots, s_n とし、それらを元とした集合 s_1, s_2, \dots, s_n を S とする
 - システム要素は、例えばインフラであればサーバ、コンテナ、ミドルウェア、ネットワークなど

ドキュメント

- システムを説明するドキュメントを集合 D とする。ドキュメントはシステム要素を説明するものである。
 - 「ドキュメント D があるシステム要素 s_n を説明できている」場合、関数 $desc()$ をつかって $desc(s_n) \in D$ とする

(関数とは集合 A の各要素に集合 B の唯一つの要素を割り当てるもの)

真に正しいドキュメント

$$desc : S \rightarrow D \quad (1)$$

また

$$desc(S) = D \quad (2)$$

(外延性の公理より、集合としては元の重複は考えない)

正しいことが記載されているが不十分なドキュメント

$$D \subseteq desc(S) \quad (3)$$

(D は $desc(S)$ の部分集合)

ある集合 A の要素の個数を記号 $n(A)$ で表すとする。

$n(D \cap desc(S))$ が $n(S)$ に近ければ近いほど良い。

$n(S) - n(D \cap desc(S)) = 0$ が理想。

間違ったことが記載されているドキュメント

$desc(s_x) \in D$ かつ $s_x \notin S$ となる元 s_x が存在する

これを本稿では「間違い要素」と呼ぶ。間違い要素の集合は、差集合 $D - desc(S)$ で表すことができる。

時間の経過に応じて変化するシステムとドキュメント

システムは継続的に開発されるものであり変化するものである。

ある時刻 t のシステムを S_t と表す。

時刻 t が、 t_1, t_2, \dots, t_n と経過する場合、それぞれのシステムを $S_{t_1}, S_{t_2}, \dots, S_{t_n}$ とすると、

簡易的に、システムの巻き戻し(revertなど)を考えないとし

$$n(S_{t_1} \cap S_{t_1}) = n(S_{t_1}) \quad (4)$$

$$n(S_{t_1} \cap S_{t_2}) > n(S_{t_1} \cap S_{t_3}) > \dots > n(S_{t_1} \cap S_{t_n}) \quad (5)$$

と、システム要素が時間の経過に応じて変化していくとする。

この場合、ドキュメント D が仮に $D = desc(S_{t_1})$ だとしても、時間の経過と共に $n(D \cap desc(S_{t_n}))$ の数は小さくなる。

これは「正しく説明できているシステム要素」の数が時間の経過に応じて小さくなることを表している。

また間違い要素集合も $n(D - desc(S_{t_n})) > 0$ として存在することになる。

システムの変化に追従するドキュメント

これは S のシステム要素の変化に応じて D の要素が追従できれば良い。

人が実施するのであれば「システムの変化に応じてドキュメントを書き直せばよい」ということになる。

ただ、近年システムの変化のスピードは早くなっている。またシステムは巨大で複雑なものになっており、手動でシステムに変化に追従するのは難しい。

そこでドキュメントの自動生成が考えられる。

まず、ドキュメント D を生成する関数を $f()$ とした場合、時刻 t に応じて変化するということは

$$D = f(t) \quad (6)$$

さらにその結果が $desc(S_{t_n})$ となれば、 S_{t_n} の変化に追従できるといえる。

$$D = f(t) = desc(S_{t_n}) \quad (7)$$

S_t の時刻 t による変化を $S_t = S(t)$ と関数で表すとする

$$D = f(t) = desc(S_{t_n}) = desc(S(t)) \quad (8)$$

$$f = desc(S) \quad (9)$$

上記よりシステム入力を入力としてドキュメントを生成できれば良いことがわかる。

システムとドキュメントを変化に応じて生成する

今までのモデルで別のアプローチも表現できる。

あるシステム S は常に設定 C から生成されるもの($gen()$)であった場合、

$$gen(C_{t_n}) = S_{t_n} \quad (10)$$

また設定 C からドキュメントも生成できれば($doc()$)、

$$doc(C_{t_n}) = D \quad (11)$$

$$C_{t_n} = doc^{-1}(D) \quad (12)$$

$$D = doc(C_{t_n}) = doc(gen^{-1}(S_{t_n})) \quad (13)$$

となり（逆関数があるとする）

$doc(gen^{-1}())$ が $f()$ といえる。

OpenAPIなどがこのアプローチにあたる

4. システムの理解のしやすさのモデル化

セクション3では、前提に「ドキュメントを通じてシステムを理解する」と置いている。

また、真に正しいドキュメントとして $D = desc(S)$ とモデル化している。

つまりシステム要素の数とドキュメントの要素の数は同数となる

$$n(D) = n(S) \quad (14)$$

近年システムは巨大で複雑なものになっており1人で全てを理解するのは現実的に難しくなっている。

つまり、一人が最終的に理解できる要素の数を m とすると

$$m < n(S) \quad (15)$$

であるということになる。

また、一人が一度に理解できる要素の数を k とすると

$$k < m < n(S) \quad (16)$$

となる。

ドキュメント D を理解しやすくするためには、部分集合 $D_1, D_2, D_3 \dots D_n$ に分割してそのそれぞれの D_n は

$$k < n(D_n) \quad (17)$$

であることが望ましいことがわかる。

部分集合の分け方については、さらに理解しやすくするために様々なアプローチが取れる。

例えばサーバー一覧やネットワーク図など、同じようなシステム要素をまとめてグルーピング化しても良いし、各マイクロサービスごとに分けても、ある処理ごとに分けても良い。

これを ISO/IEC/IEEE 42010:2011 にならって「ビュー」と呼ぶ（ただし本家のそれとは厳密には定義が異なる）。

ビューはステークホルダー（そのドキュメントの対象者）が理解すべきコンテキストに沿って記述される。

ネットワーク管理者にとってはネットワーク図、あるマイクロサービスを管理するSREにとってはそのマイクロサービスを構成するノードとネットワーク、など。

よって D_1 に含まれる要素が D_2 に含まれることもある。

よって理解しやすさを持ったドキュメントの特徴としては以下ようになる

$$D = D_1 \cup D_2 \cup D_3 \cup \dots \cup D_n \quad (18)$$

5. ドキュメントシステムに必要な要件

「継続的にシステムの変化に追従できる」「理解しやすい」ドキュメントシステムの要件として

セクション3より、

$$D = f(t) = desc(S_{t_n}) = desc(S(t)) \quad (19)$$

であること、セクション4より、

$$D = D_1 \cup D_2 \cup D_3 \cup \dots \cup D_n \quad (20)$$

というドキュメントが生成できることと主張した。

ここで、さらに $f(t)$ について $desc(S)$ とするために分解すると、

1. 入力 S からシステム要素を列挙すること
2. 1で列挙したそれぞれのシステム要素について理解のための文章を追加すること（本稿ではアノテーションと呼ぶ）

と分解できそうである。

ここで、2のアノテーションは技術的に困難であるとして人に任せるとすると、1だけが要件として残る。

つまり、

1. 対象システムからシステム要素を抽出でき、それぞれに対して穴埋め形式でドキュメントを記述できる
2. 1で抽出したシステム要素を任意のグループでグルーピングすることができ、それに対してドキュメントを記述できる
 - システム要素はグループ間で重複可
 - 1で抽出した要素のみで2を実現することで $s \notin S$ となる要素が存在しないように

という要件が見えてくる。

ndiag

上記の要件を満たすドキュメンテーションツールとして筆者はndiagというツールを作成している

<https://github.com/k1LoW/ndiag>

ndiagはシステム要素の抽出機能を持っていないので、例えばノード間のリレーションなどはサービスマップや<https://github.com/yuuki/shawk>などを入力とするようなことを想定している

参考資料

- [目的に沿ったDocumentation as Codeをいかにして実現していくか / PHPerKaigi 2021](#)
- [ソフトウェアシステムアーキテクチャ構築の原理 第2版](#)