

Universal SAX : 空間充填曲線を使用した SAX の多次元時系列データへの適用

理学専攻・情報科学コース 大西 史花

1 はじめに

センサーデバイスやシミュレーション技術の発展に伴い、GPS ログなどの軌跡データ、車の動きなどのムービングデータなど大量の多次元時系列データの入手が容易になっているが、これら膨大なデータの活用には効果的な問い合わせや索引が不可欠である。その中の 1 つに Lin らが提案した Symbolic Aggregate Approximation(SAX)[2] がある。SAX は時系列データを文字列に変換する。この手法はアルゴリズムが単純である一方で元データに対してマイニングを行った結果と同等の結果を生成することが出来るなど有用な特徴を多々持つ。しかしながら、SAX は 1 次元の時系列データしか扱うことができない。既存の手法では多次元時系列データの次元を重心距離法や主成分分析を用いて 1 次元に次元削減した後 SAX を適用しているが [3][4]、これらの手法はデータ間の位置関係を損ないマイニングに必要なデータを欠落させてしまう場合がある。

よって本稿では、空間充填曲線を用いて多次元空間内のデータの位置関係を損なうことなく文字列に変換する Universal SAX (USAX) を提案する。空間充填曲線とは、多次元空間を曲線で埋め尽くすように辿ることで空間の順序付けを行う手法である。本手法において文字列に変換後の各文字は元のデータ空間に適用した空間充填曲線による順序値の範囲によって割り当てられる。空間充填曲線による順序値の範囲は元空間において多面体を成すため、2 つの USAX 文字間の距離は 2 つの多面体間の最小距離として計算することが出来る。このように、USAX 文字列化後の距離測度として元の多次元空間上の距離を用いるので、本提案手法による多次元時系列データの 1 次元化は元データのマイニング等に有用な情報を欠落させることなく行うことが出来る。

2 SAX

SAX とは、時系列データの表現手法の 1 つで、時系列データを文字列に量子化する手法である。SAX による時系列データの文字列化を図 1 によって示す。図 1 では破線が SAX を適用する時系列データである。SAX を使用して時系列データを文字列に変換する手順は以下の通りである。まず時間軸を等間隔に区分し(縦実線)、区間毎の平均値を算出する(横太線)。その後、正規分布の各面積が等しくなるような分割線(横実線)を定め、その区間に、a, b, c...とアルファベットを割り振る。最後に、求めた平均値を該当する区間に割り振られている文字に変換する。なお、時系列データは位置やスケールのズレを修正するためにあらかじめ正規化されている。

以上のような手順で変換される SAX 文字列は、本来実数値である時系列データに対して文字列用の検索・分析手法が適用できるという特徴を持つ。また、アルゴリズムが単純である一方で元データに対してマイニ

ングを行った結果と同等の結果を生成することが出来る。しかしながら、SAX は 1 次元の時系列データしか扱うことができない。

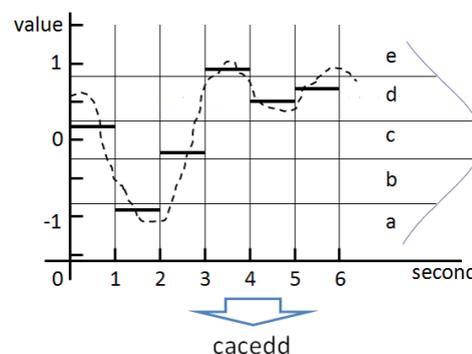


図 1: SAX による 1 次元時系列データの文字列化

3 Universal SAX

本節では多次元時系列データを文字列に変換する手法である Universal SAX を提案する。本提案手法では空間充填曲線という手法のひとつであるヒルベルト曲線を用いている [1]。空間充填曲線は多次元空間から 1 次元空間へのマッピング方法として広く使われている手法である。USAX 文字列の作成手順の概略を 2 次元データを例にして図 2 に示す。まず対象データをそれぞれの次元ごとに標準正規分布に正規化し、対象データ(黒実線)を t_0, t_1, \dots, t_n まで等間隔に区分、区間ごとの中央値 $x_0, x_1, \dots, x_{(n-1)}$ (黒点)を決定する(図 2 (1))。次に多次元空間を $2^k \times 2^k$ の格子状になるよう量子化する(図 2 (2))。その後ヒルベルト曲線を用いて正規化された対象データを区間に対応する数列(この例では $\{51, 46, 40, 38, 27, 28, 10, 17, 15, 1\}$)に変換、時系列状のデータとして表す(図 2 (3))。最後に、各次元で行った正規化の領域が等しくなるよう分割を行い、それぞれの領域にアルファベットを割り振り(図 2 の(4)左図)文字列に変換する(図 2 の(4)左図下)。図 2 の(4)の右図は正規分布を分割した後の多次元空間のイメージである。

Universal SAX 文字列間の距離計算は、その文字列を構成する文字間の距離の合計として行われる。文字間の距離は、それぞれの文字が割り振られている空間上の領域間の最小距離として計算される。例えば図 2 の(4)の右図のように分割された空間上での文字 a と d の距離 $dist("a", "d")$ は、マス目 1 つを距離 1 とすると 1 となり、文字 a と g の距離 $dist("a", "g")$ は空間上で領域が隣接しているため 0 となる。このように、USAX は元の空間上の距離関係を用いて文字列データに変換しているため文字列変換時にデータ間位置情報の欠落が起きない。

