

球状蛋白質分子の構造量におけるカオスの性質の検出

(Detection of Chaotic Aspects in Structural Quantities for Globular Proteins)

○ 窪田 綏 (Advanced Tech. Inst., Inc.)、土屋 尚 (明星大・情報学部)

〔要旨〕

Previously we showed [Y.Kubota, et al., *Biphsica. Acta*, Vol.1079, 73-78 (1991)] the distance ρ of each $C\alpha$ atom from the center of mass and the quantity called N14 [K.Nishikawa and T.Ooi, *J. Biochem.* Vol.100,1043-1047 (1986)], which represents the number of all the $C\alpha$ atoms within a sphere of 14\AA radius centered on each $C\alpha$ atom, are the structural quantities which reflects the tertiary structure of a protein molecule from its amino acid sequence by the method proposed by Grassberger and Procaccia (the G-P method) [P.Grassberger and I.Procaccia, *Physica D*, Vol.9, 189-208 (1983)] in the chaos theory. That is, the correlation dimension D_2 and the embedding dimension d become 4.7 and 20, respectively for the sequences of ρ and N14 of globular proteins. According to our remark, the experimentally obtained sequence of the quantities has been generated by a certain deterministic mechanism. In other words, it is suggested to be deterministic chaos.

Here, we try another approach than the above G-P method, in order to detect a chaotic behavior of protein structures. Briefly, we construct a predictor from the data of K residues long from the sequence of the structural quantity by using RBF (Radial Basis Function) method proposed by Casdagli [M. Casdagli, *Physica D* Vol.35, 335-356 (1989)]. By applying this predictor, we can find the short-term predictability and the long-term unpredictability in globular proteins, which is a characteristic feature of the deterministic chaos.

Furthermore, we can generate a data over the number of residues of a protein by iterate to use the predictor, recursively (free running). Interestingly, the data which is obtained in such a way shows a similar attractor to the original one. And also the application of the G-P method to it reveals the same the correlation dimension D_2 and the embedding dimension d to the native ones of a protein.

These results suggest the existence of chaotic aspects in globular proteins.

〔序論〕

1. 以前我々〔1〕は、球状蛋白質の残基がどの程度表面―内側に位置するかの度合いを示す指標として各々 C^α 原子を中心とした半径 1.4 \AA の球内に含まれる全ての C^α 原子の数で定義される量 N_{14} 〔2〕や重心からの各 C^α 原子の距離 ρ がアミノ酸配列に沿った量の（1次元）数列として、物理化学的な決定論的メカニズムによって生成されると考えられている立体構造を反映するStructural Quantityであることをカオス理論の手法であるGrassberger と Procacciaによる相関積分の方法(G-P法)〔3〕により示した。つまり球状蛋白質における N_{14} や ρ 等の量の列はランダムというよりカオスの振舞いを見せることを示した。従って我々は蛋白質の構造形成にはある決定論のプロセスによって制御されているに違いないと考える。
2. 今回この考えを更に裏付けるために、我々は別な方法である Casdagli の動径基底函数 (Radial Basis Function: RBF) による手法〔4〕を使ってこれらの量にカオスの性質が存在することを示すことが出来た。即ち、RBFによりこれらの量を Time Series と見做して予測を試みた時、短期予測は可能であるが長期予測は不可能という意味においてカオスの性質があることがわかった。
3. また残基数 N の蛋白質においてある学習データから更にその先 N 残基分フリーランさせてそのアトラクターを求めたところそのオリジナルアトラクターと定性的に究めて類似のアトラクターを得ることができ、更にこれを定量的にG-P法による相関積分によってCorrelation Exponentを求めそこから埋め込み次元 d と相関次元 D_2 を求めたところオリジナル蛋白質の場合とほぼ同じプロフィールと次元が得られ、典型的なカオス系の場合と同様にストレンジアトラクタの再構築ができた
4. このことから N_{14} や重心からの各 C^α 原子の距離 ρ のような構造量にカオスの性質が存在することが分かり、益々決定論的なメカニズムで蛋白質構造が形成されていることが支持されたと考える。

〔方法〕

○ Casdagli による非線形予測法

Casdagli は動径基底函数 (Radial Basis Function: RBF) 法をカオスの時系列の予測問題に応用した。即ち;

$v(t)$ を再構成されたアトラクターの軌道の 1 点, c_i をセンターと呼ばれる点、とした時 1 期先の予測を

$$\hat{y}(t+1) = \hat{F}(v(t)) = \sum_i^M \tilde{\gamma}_i \Phi(\|\tilde{v}(t) - \tilde{c}_i\|) \quad (1)$$

により表す。ここで、 γ_i は重み、 $\Phi(\|\tilde{v}(t) - \tilde{c}_i\|)$ は基底函数と呼ばれる。基底函数としては;

ガウシアン函数: $\Phi(r) = \exp(-r^2 / b) \quad (2)$

ベル型函数: $\Phi(r) = \frac{1}{1 + \cosh(\sigma_h r)} \quad (3)$

が用いられる。

さて、通常データは 1 次元の離散的データとして得られるので、Eqn.1 を行列・ベクトルの形で表現すると便利である。即ち、各ベクトル \tilde{x} を m 次元ベクトルして;

$$\mathbf{y} = (\tilde{y}_2, \dots, \tilde{y}_{M+1})^t \quad (4)$$

$$\boldsymbol{\gamma} = (\tilde{\gamma}_1, \dots, \tilde{\gamma}_M)^t \quad (5)$$

$$\mathbf{P} = [p_{ij}] = \Phi(\|\tilde{v}(i) - \tilde{c}(j)\|) \quad (6)$$

とすると

$$\mathbf{y} = \mathbf{P}\boldsymbol{\gamma} \quad (7)$$

となり、予測器の作成は Eqn.7 の $\boldsymbol{\gamma}$ を求め、これを次々に再帰的に繰り返し用いることで長期に渡って時間発展 (Free Running) させて、予測時系列を作成することができる。こうして、与えられた時系列の決定論的力学系としての挙動の特徴を捉えることができる。なお、今後予測器 (predictor) とは Eqn.1 または Eqn. 7 を指して言う。

○ 相関積分法

Grassberger と Procaccia らは、解析対象となるシステムの有する自由度であるアトラクターの次元を推定する相関積分の手法を提案した [3]。即ち、今 1 次元時系列を $\{x(t_i)\}$ とし、この 1 次元データから m 次元ベクトルを;

$$\vec{X}_i = (x(t_i), x(t_i + \tau), \dots, x(t_i + (m-1)\tau)) \quad (8)$$

を構成する。ここで τ は時間ラグである。この時、相関積分 $C(r)$ を；

$$C^m(r) = \frac{1}{N^2} \sum_{i,j=1}^N \Theta(r - |X_i - X_j|) \quad (9)$$

と定義する。ここで、 N はデータ数、 Θ はヘビサイド関数である。この $C^m(r)$ が；

$$C^m(r) \propto r^{\nu(m)} \quad (10)$$

と表わされる時、 $\nu(m)$ を相関指数と呼ぶ。この $\nu(m)$ は Eqn.10 の両辺の対数；

$$\log C^m(r) \propto \nu(m) \log r \quad (11)$$

をとり、横軸に $\log r$ 、縦軸に $\log C^m(r)$ をプロットした時、適切な r の範囲内での直線部分の勾配として求められる。そして相関次元 D_2 は、空間次元 m を漸次増加させて行った時、もはや $\nu(m)$ が増加しない縦軸の点 ν として求められる。また ν がこれ以上変化しない最小の m は埋め込み次元 d とよばれる。

〔結果〕

Fig.1 は典型的な球状蛋白質の例として porcine elastase (PDB code: 1est 残基数 $N=240$) に対して得られた西川及び大井によるパラメータ N_{14} のプロット [2] である。

さて、Fig.1 に示された 1 次元時系列に前項で導入した Casdagli の予測法を適用して見る。即ち、Eqn.1 における学習データ長 M を 194 とし、また前回我々によって得られた結果を参考にして再構成されたアトラクターの次元 m の値を 6 として [1]、予測器、Eqn. 1 を Free Running させた結果を Fig.2 に示す。残基番号 201 (= $M+m+1$) 以降がフリーラン予測となるが、残基 201 ~ 220 付近の予測精度は高く、それ以降となるとその予測精度は漸減しており、カオスとしての特徴である短期予測可能性と長期予測不能性と言うカオスの挙動を良く呈していると考えられる。

さて、この予測モデルがオリジナルな N_{14} の挙動を良く再現しているかどうか、まず得られたアトラクターの様子から定性的に見てみよう。Fig.3a はオリジナルなアトラクターであり、Fig.3b はフリーランによるそれである。この予測モデルが比較的良くオリジナルの振る舞いを再現していることが分かる。更に定量的に見るために G-P 法による相関積分を計算することにより、相関次元 D_2 及び埋め込み次元 d を求めて見る。

Fig. 4 a はオリジナル N_{14} (Fig.1) 、 Fig. 4 b は 201 残基からフリーランさせたものについて (Fig.2) $m=1\sim 30$ に対する相関積分 (Eqn. 9) である。

Fig. 4 a (オリジナル) 及び Fig. 4 b (フリーラン) からそれぞれの相関次元 D_2 及び埋め込み次元 d を求めるために、これらの図のカーブの直線部分の勾配 $v(m)$ を計った (Fig. 5)。

Fig. 5 から示されるようにオリジナル (点線)、フリーラン (実線) 共に相関次元 D_2 及び埋め込み次元 d がそれぞれ 5 と 20 となっており、この予測モデルが良くそのオリジナルな振る舞いの特徴を再現していると言える。また重心からの各 C^α 原子の距離 ρ についても類似の結果が得られた。

以上のことから N_{14} や重心からの各 C^α 原子の距離 ρ のような構造量にカオス的性質が存在すると考えられ、決定論的なメカニズムで構造が形成されていると考える。このことは、蛋白質の構造形成過程の背後に、恐らくは自由度が 20 であるような微分方程式或いは差分方程式の存在を暗示していると思われる。

[参考文献]

- [1] Kubota, Y. and Tsuchiya, T. : "Chaos-Theoretical Analysis of Possible Structural Quantities for Globular Proteins," *Biophysica Acta*, **1079**, 73-78 (1991)
- [2] Nishikawa, K. and Ooi, T. : "Radial Locations of Amino Acid Residues in a Globular Proteins: Correlation with the Sequence" *J. Biochem*, **100**, 1043-1047 (1986)
- [3] Grassberger, P. and Procaccia, I.: "Measuring the Strangeness of Strange Attractors" *Physica D* **9**, 189-208 (1983)
- [4] Casdagli, M. : "Nonlinear Prediction of Chaotic Time Series", *Physica D* **35**, 335-356 (1989)

Figure Legends

Fig.1 Porcine Elastase (残基数240、PDB Code: 1est) の N_{14} プロフィール。各残基に対してPDB座標データから14 Å以内の C^α 原子がカウントされている。表面に近い残基ほど N_{14} の値は高く、内側のそれは低く、丁度各残基の C^α 原子の重心からの距離 ρ とほぼ逆相関の関係になっている。

Fig. 2 Fig. 1において、残基番号1~200のデータを使って予測器 Eqn.1 を作成し、201番目の残基からもう240先までフリーランさせた。実線はフリーラン、破線はオリジナル。201~220番目付近までは良く予測されており（短期予測可能性）、それ以降の予測はずれており（長期予測不能性）、決定論的カオスである特徴を呈している。

Fig. 3 a-b (a) N_{14} プロフィール (Fig. 1) から得られたオリジナルなアトラクター、(b) フリーラン (Fig. 2) から得られたアトラクター。定性的には、オリジナルアトラクターの振る舞いを再現する予測モデルが構築できていることが分かる。

Fig. 4 a-b 再構成空間次元 m を1から30まで変化させた時の相関積分 Eqn 9。
(a) オリジナル (b) フリーラン

Fig. 5 各 m に対して Fig. 4 a, 4 b の直線部分の勾配から得られた相関指数 $\nu(m)$ のプロット。実線はフリーラン、点線はオリジナル。このグラフから共に、相関次元 $D_2 = 5$ 、埋め込み次元 $d \cong 20$ と求められた。

Fig.1

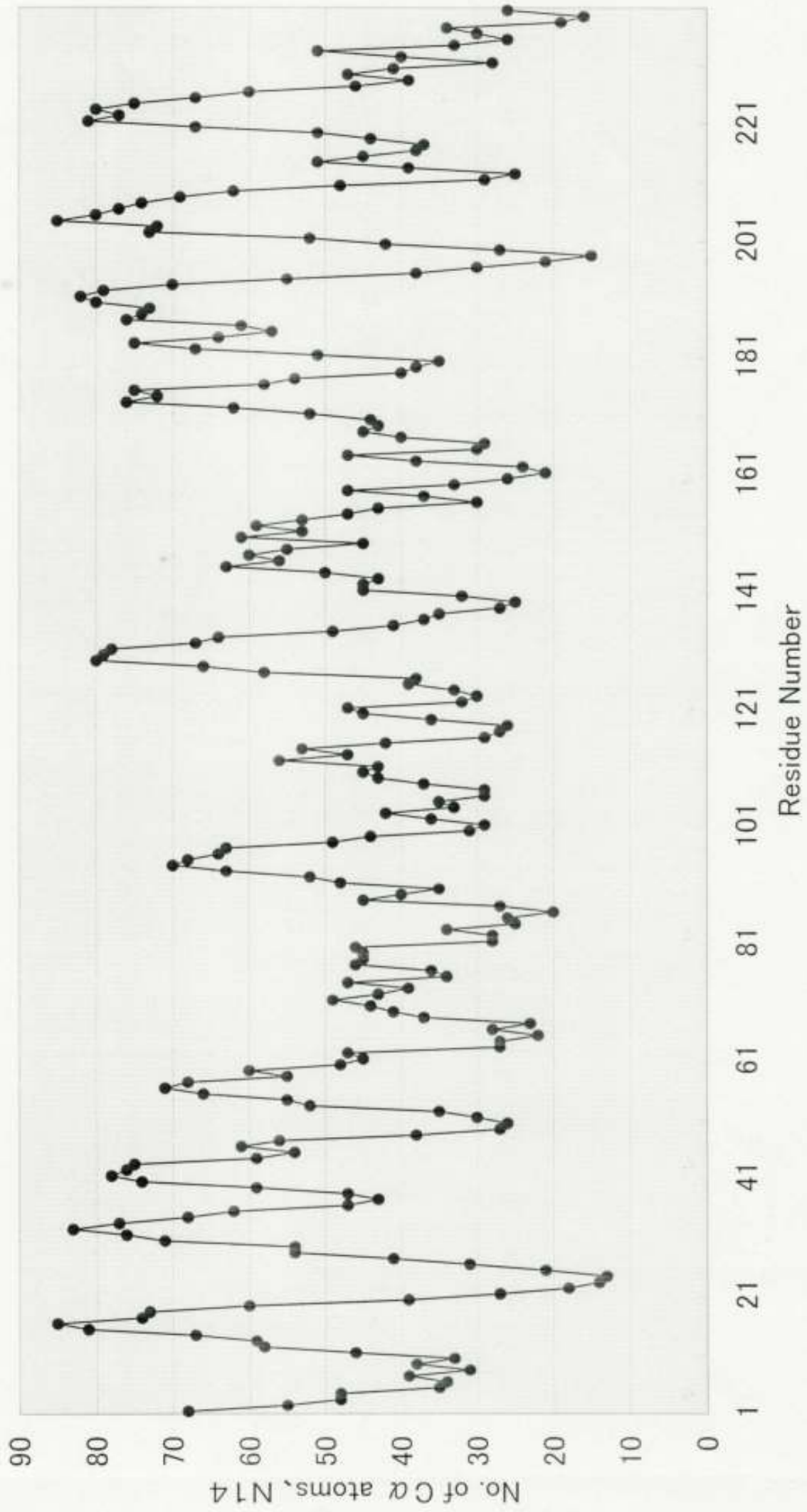


Fig.2

—●— Free Run
- - -▲- - - Original

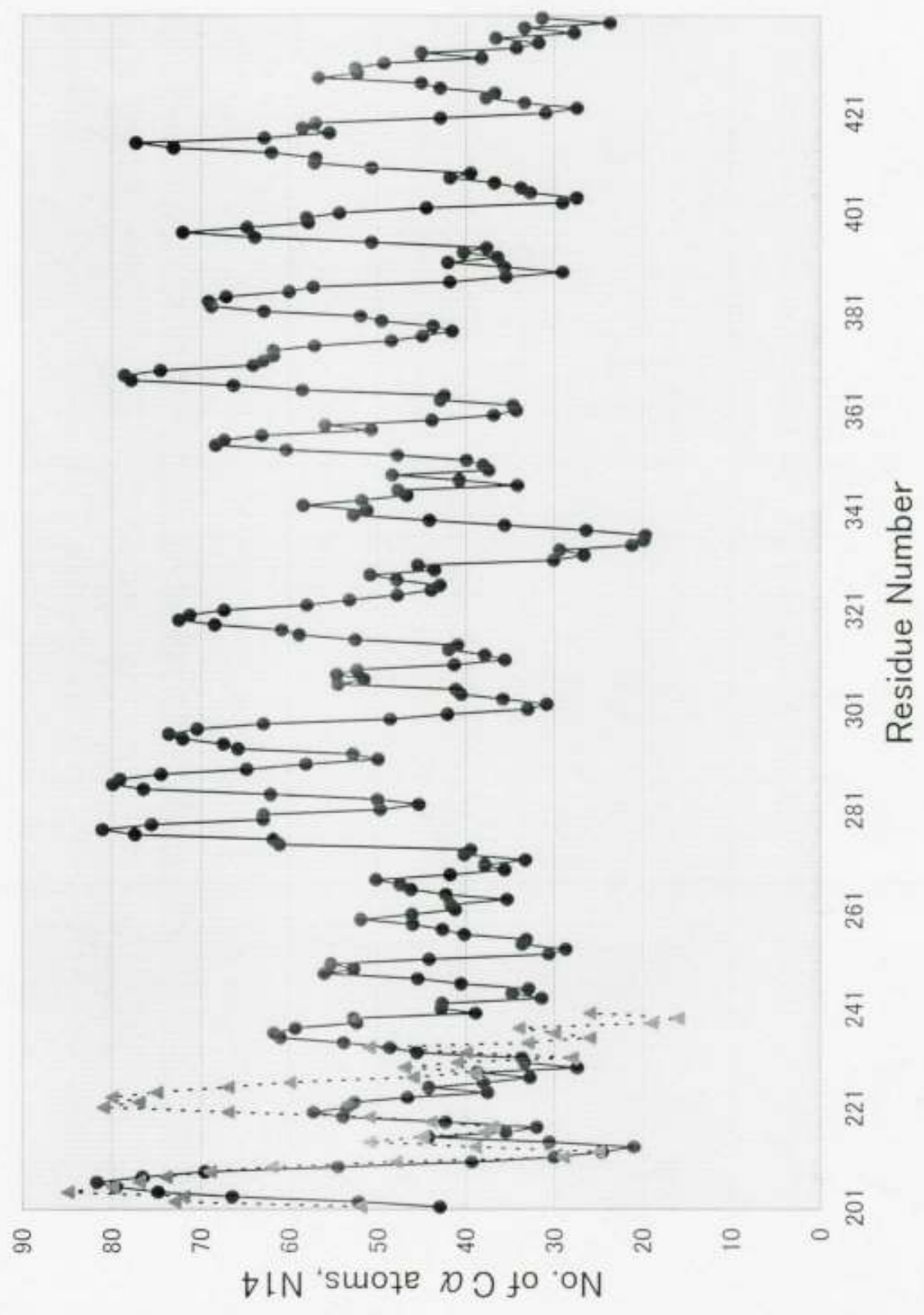


Fig.3a (Original)

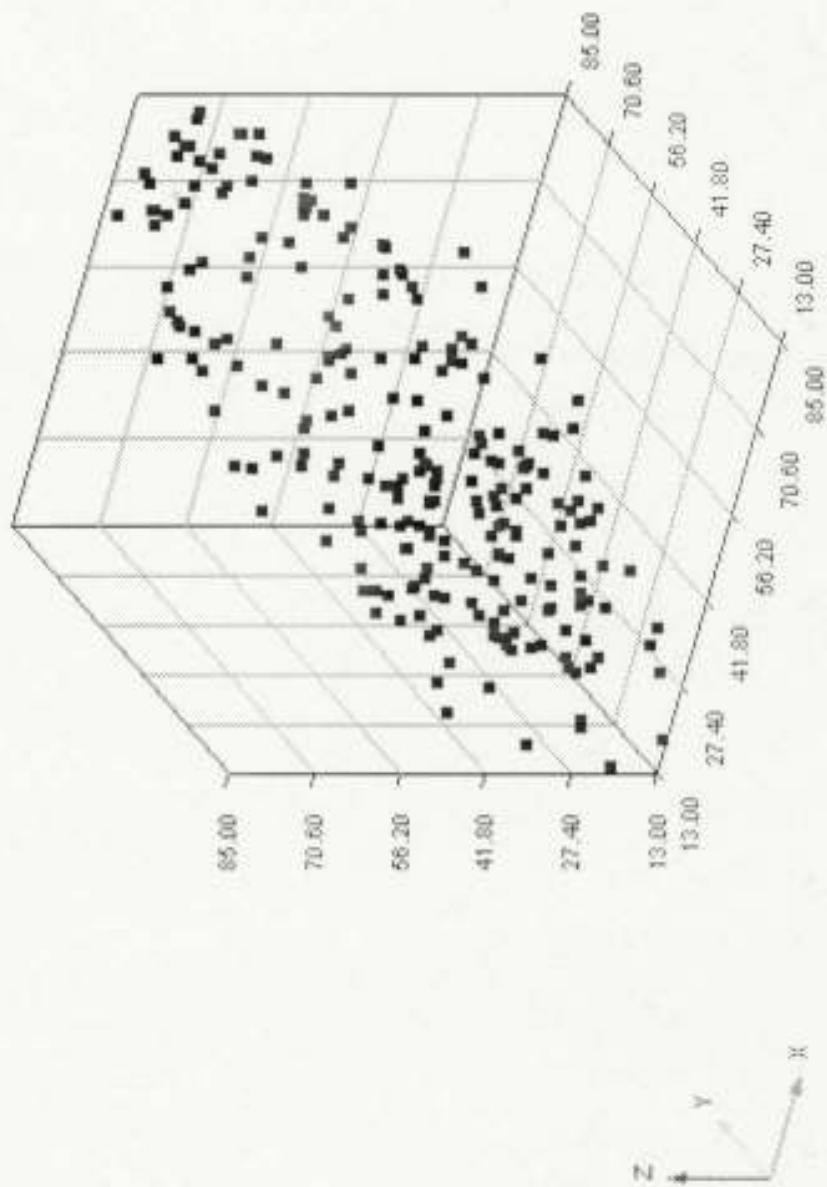


Fig.3b (Free Run)

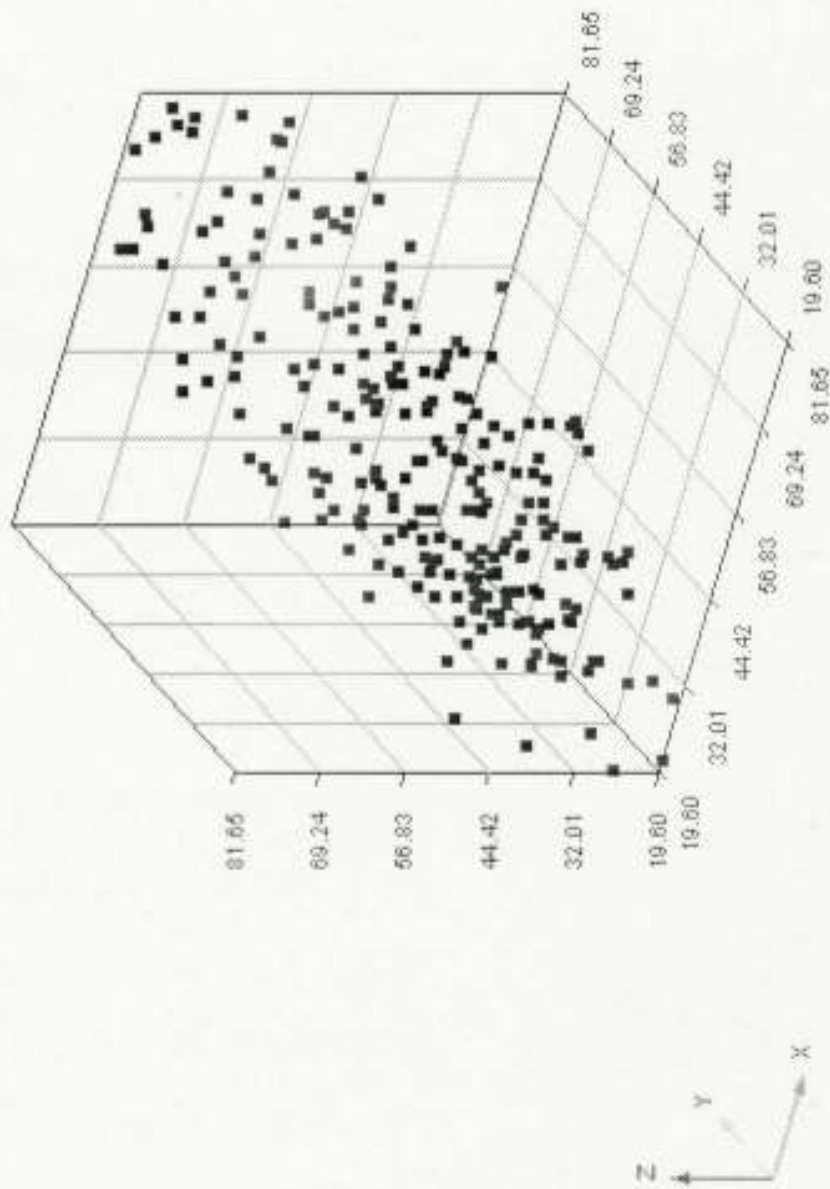


Fig. 4a

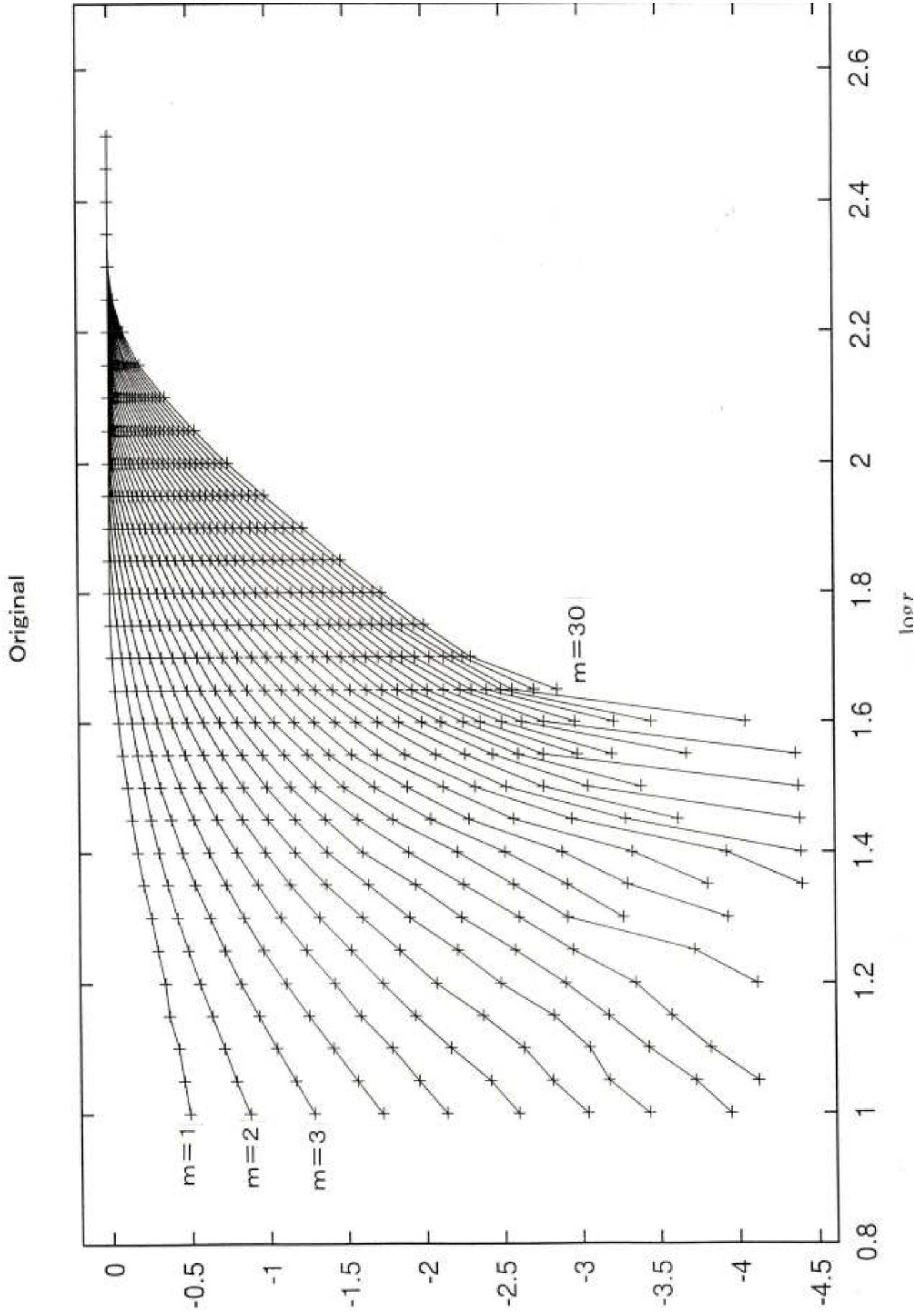


Fig. 4b

Free Run

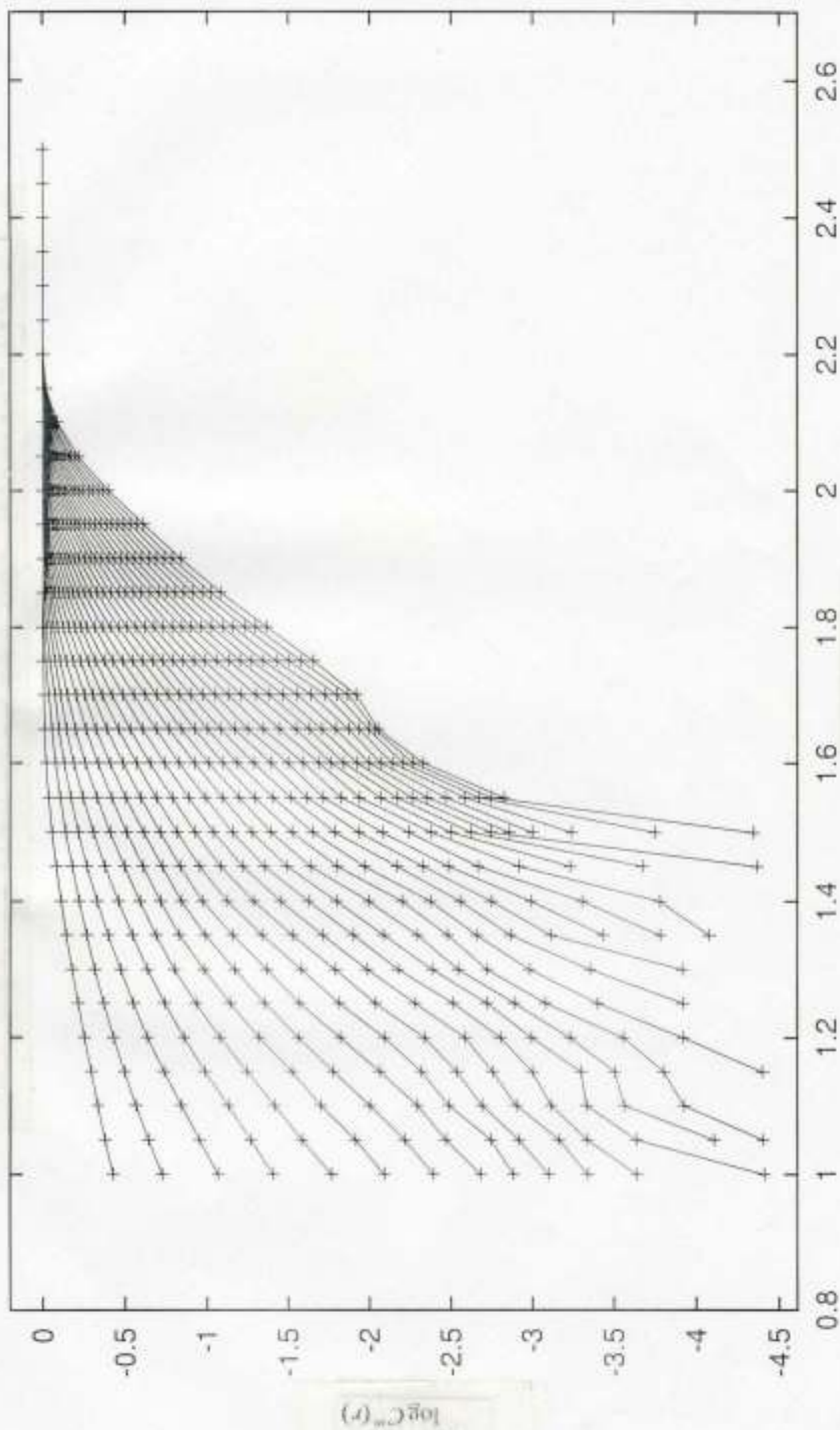


Fig. 5

