

線形回帰法を用いたタンパク質原子電荷に関する研究

(東大院工*、東大生研**) ○金泰煥*、平野敏行**、佐藤文俊**

Studies of atomic charges for proteins using the linear regression method

(School of Engineering, the University of Tokyo*,
Institute of Industrial Science, the University of Tokyo**)
OKim Taehwan*, Hirano Toshiyuki**, Sato Fumitoshi**

【序】 ESP 電荷[1]は分子周辺の静電ポテンシャル(ESP)を再現するよう、フィッティングした原子電荷である。評価点*i*において、カノニカル分子軌道(CMO)計算結果から得られた ESP を V_i 、原子電荷 q から算出した ESP を \hat{V}_i とする。ESP 電荷は最小二乗法より残差 $(V_i - \hat{V}_i)$ の2乗和である l_2 損失 $f_{LS}(q)$ が最小となる原子電荷である(式(1))。

$$f_{LS}(q) = \sum_{i=1}^{N_{grid}} (V_i - \hat{V}_i)^2 \quad (1)$$

ESP 電荷の評価点*i*は全て vdW 半径の外側に分布するため、分子内部は特異な電荷が得られる場合がある。このような最小二乗法の欠点を防ぐために、様々な線形回帰法が開発されている。本研究では、線形回帰法を用いてタンパク質の原子電荷を算出し、特性を調べた。

【理論】 l_2 制約付き最小二乗学習法である Ridge 回帰は l_2 損失に、 l_2 ノルムの正則化項 $\|q\|^2$ を加え、過適合を防いだ回帰法である(式(2))。Ridge 回帰の l_2 ノルムの代わりに l_1 ノルムの正則化項 $\|q\|_1$ を用いた l_1 制約付き最小二乗学習法(Lasso 回帰法)は多くのパラメータ q がスパースと推定される学習法である。パラメータ q の次元が大きいとき、効率良く特徴選択ができる。式(2)、(3)の λ_R 、 λ_L はそれぞれ Ridge パラメータ、Lasso パラメータとよび、損失項と正則化項のバランスを表す。

$$f_R(q) = f_{LS}(q) + \lambda_R \|q\|^2 = \sum_{i=1}^{N_{grid}} (V_i - \hat{V}_i)^2 + \lambda_R \sum_{A=1}^{N_{atom}} q_A^2 \quad (2)$$

$$f_L(q) = f_{LS}(q) + \lambda_L \|q\|_1 = \sum_{i=1}^{N_{grid}} (V_i - \hat{V}_i)^2 + \lambda_L \sum_{A=1}^{N_{atom}} |q_A| \quad (3)$$

Huber 損失 $\rho_H(r_i)$ を用いたロバスト回帰は、残差 $r_i = (V_i - \hat{V}_i)$ の絶対値が閾値 η の以下であれば l_2 損失を、大きければ l_1 損失を使う Huber 損失 $\rho_H(r_i)$ の総和が最小とする回帰法である(式(4))。 l_1 損失を混ぜ合わせたので高いロバスト性を持ち、閾値 η を設定することで、モデルの再現性とロバスト性のバランスを調整することができる。

$$f_H(q) = \sum_{i=1}^{N_{grid}} \rho_H(r_i) \quad \rho_H(r_i) = \begin{cases} \frac{r_i^2}{2} & \text{if } |r_i| \leq \eta \\ \eta|r_i| - \frac{\eta^2}{2} & \text{if } |r_i| > \eta \end{cases} \quad (4)$$

【方法】 本研究はオキシトシン(原子数:134 個、評価点数:5,414 個)とインスリン(原子数:786 個、評価点数:23,787 個)を計算対象にした。CMO 計算および ESP 計算は ProteinDF[2]を使用した。Ridge

回帰、Lasso 回帰、Huber 損失 $\rho_H(r_i)$ を用いたロバスト回帰から見積もられた電荷を便宜上 Ridge 電荷、Lasso 電荷、Huber 電荷とする。Ridge 電荷は制限関数を調和型、参照電荷を 0 にした RESP 電荷[3]に相当する。全ての原子電荷計算は原子電荷の総和が系全体の電荷 Q_{tot} に等しい条件下 $g(q)$ (式(5))、ラグランジュ未定乗数方程式から原子電荷 q を求めた。

$$g(q) = \sum_{A=1}^{N_{atom}} q_A - Q_{tot} = 0 \quad (5)$$

線形回帰法から求められた原子電荷を Mulliken 電荷および Amber 力場電荷 ff03 電荷と比較ならびに解析を行った。

【結果】 インスリンの Mulliken 電荷、ff03、Ridge 電荷 ($\lambda_R = 5 \times 10^{-4}$) の計算結果を図 1 に示した。ここでは主鎖の窒素(N)と酸素(O)のみ示した。Mulliken 電荷において N と O の原子電荷は同程度であり、分散が小さいのに対し、Ridge 電荷は隣接する原子が作る ESP によって電荷が見積もられるため、多様な原子電荷を取ることができる。51 残基すべての主鎖の N と O の平均値を調べると N が O より正に寄っていることが確認された。このことは、N と O の電気陰性度から考えると妥当な結果である。

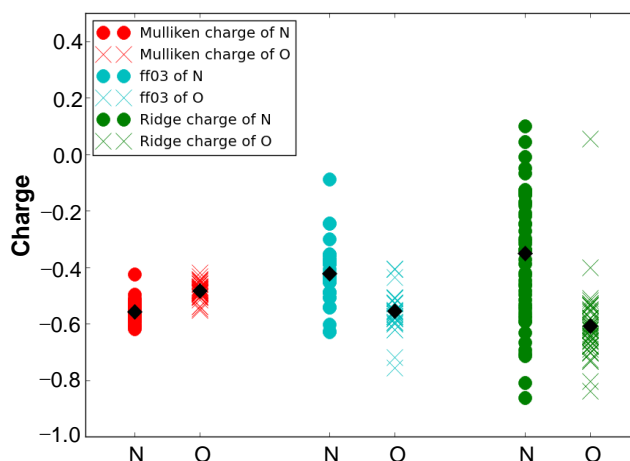


図 1 インスリンの主鎖の窒素と酸素の電荷
(黒いダイヤモンドは電荷の平均値)
(左)Mulliken 電荷, (中)ff03, (右)Ridge 電荷

次に様々な Lasso パラメータ λ_L を使って Lasso 電荷がスパースに推定されたパラメータ数 (0 要素数) と ESP 再現性の指標となる RRMS の結果を表 1 に示した。Lasso パラメータ $\lambda_L = 10^{-1}$ のとき、全体の 75% に相当する 596 個の原子がスパースに推定されたにも関わらず、ff03 と同程度の ESP 再現性を見せた。即ち、Lasso 電荷は少ない変数 (原子電荷) で ff03 と同程度の ESP が再現できるといえる。

Huber 電荷の特徴については当日報告する。

表 1 様々な Lasso パラメータ λ_L を使って
Lasso 電荷がスパースに推定されたパラメータ数 (0 要素数) と ff03 の RRMS

原子電荷	Lasso 電荷					ff03
	$\lambda_L = 0$	$\lambda_L = 10^{-4}$	$\lambda_L = 10^{-3}$	$\lambda_L = 10^{-2}$	$\lambda_L = 10^{-1}$	
0 要素数	0	41	202	372	596	0
RRMS	0.0284	0.0287	0.0340	0.0520	0.1168	0.1210

$$RRMS = \sqrt{\frac{\sum_{i=1}^{N_{grid}} (V_i - \hat{V}_i)^2}{N_{grid}}} \quad (6)$$

【参考文献】 [1] U. C. Singh. et. al., *J. Comput. Chem.*, 5 (1984) 129. [2] <http://proteindf.github.io/>. [3] C. Bayly. et. al., *J. Phys. Chem.*, 97 (1993) 10269.