DNNを用いて抽出された力学系の対称性からの保存量推定

Conservation-Law Estimation from Symmetric Property of Dynamical System Extracted using DNN

本武 陽一

Mototake Yoh-ichi

東京大学大学院新領域創成科学研究科

Graduate School of Arts and Science, The University of Tokyo

It is suggested that Deep Neural Networks (DNN), which continues to develop in recent years, has a function to extract information of data sets necessary to achieve a given task by modeling the distribution as a manifold. In addition to confirming the usefulness of DNN technology, numerous researchers and engineers are developing various DNN algorithms and tuning parameters. This situation means that enormous knowledge on the manifold structure for various data sets is being accumulated. The purpose of this research is to propose a method to extract manifold structure with complex shape extracted in an interpretable form. Specifically, we propose a method to extract the symmetry of manifold for coordinate transformation. Based on the Noether's theorem in physics, we develop the method to estimate the conservation law of the system. Applying the proposed method to the time series data of the moving object according to the central force potential, it was confirmed that symmetry according to the conservation law of angular momentum could be extracted.

1. 導入

近年発達を続ける Deep Neural Networks(以下, DNN) が、与えられたタスクを達成するために必要なデータセットの 情報を、その分布を多様体としてモデル化することで抽出す る機能を持つことが示唆されている [1][2][3][4][5][6][7][8]. こ のような視点に立つと、DNN の写像関数は多様体を定義する ものと理解される.しかしながら、この写像関数は一般に、非 常に多数のパラメータによって構成された非線形関数となるた め、そこからデータセット分布の持つ情報を解釈可能な形で抽 出することは困難である.

これに対して我々は、この DNN がモデル化した多様体構造 についての情報を、その入力空間の連続的な変換に対する対称 性として抽出する手法を提案する.これによって、複雑な写像 関数を対称性という指標でより単純に評価することが可能と なる.

対称性から対象のデータセットの理解を行うためには,対称 性とデータセットの性質との間を結ぶ理論が必要となる.物理 学の分野では,この対称性と対象の性質に関する理論の構築が 長年行われて来た.ネーターの定理はその代表例で,ハミルト ニアンの座標変換に対する普遍性と,系の保存量とを結びつ けることが可能となる.本研究では,提案手法を中心力ポテン シャル系の時系列データに適用し,そこからハミルトニアンの 対称性を抽出することを試みた.この結果,角運動量に対応す る対称性が抽出されることを示唆する結果が得られた.

2. 手法

2.1 DNN と多様体

多様体は、ユークリッド空間を貼り合わせることによって構成される空間である。大まかな例としては、地図と二次元の多様体である地球表面がある. 我々は、地球表面を2次元のユークリッド空間である地図のラミネーションとして理解する.より正確には、dm 次元微分可能多様体 M とは、多様体上の任

意の点 **x**⁰ で接空間 T_x₀ M と呼ばれる d 次元ユークリッド空間と同相な近傍を持つ位相空間である.本研究では,DNN は データセットの分布を,多様体として近似的にモデル化すると 考える.その上で,DNN はデータ多様体をそれと同じ次元の ユークリッド空間に近しい大域的座標系へ写像する機能を持つ と考える.本研究では,データ分布を多様体として近似した際 に現れる構造をデータ多様体と呼ぶ.

微分同相写像は、d_m 次元多様体 M を d_m 次元多様体 N へ 写像する微分写像関数として定義される. 逆関数定理より, 微 分写像関数のうち,そのヤコビアン行列が多様体 M 上の全て の点で正則行列となるものは微分同相写像となる. つまり, あ る微分写像関数のヤコビアン行列が入力空間と出力空間上の 全ての点で正則行列となる場合、入力空間と出力空間は多様体 としてモデル化されているということができる. 多様体仮説 で取り扱うのは、ユークリッド空間中の多様体となる.この場 合,多様体を定義する微分写像に求められる条件は,多様体に 対応する部分空間から多様体に対応する部分空間への写像に 関して同相写像となることである.これは前の議論と同様に, ある微分写像が与えられた場合に、その部分空間上の全ての 点に関する写像関数のヤコビアンが正則行列となることと同 値となる。つまり、ある部分空間上の全ての点での写像関数の ヤコビアン行列の rank が dm となる場合,その部分空間が多 様体としてモデル化されているといえる. これに基づき, 我々 は DNN 写像関数がモデル化した入出力空間での多様体の接 空間を以下の手続きで算出した. 今、入力が din 次元, 出力が dout 次元である, L 層 DNN を考える. この DNN の写像関 数 $\mathbf{F}(\mathbf{x}) = (F_1(\mathbf{x}), F_2(\mathbf{x}), \cdots, F_{d_{out}}(\mathbf{x}))$ は、

$$F(\mathbf{x}) = h^{L} = f(w^{L-1}h^{L-1} + b^{L-1})$$

= $f(w^{L-1}f(\cdots f(w^{1}x + b^{1})\cdots) + b^{l-1})$ (1)

によって与えられる.ここで、 $h^l = (h^l_0, h^l_1, \cdots, h^l_{d_l})$ は、 d_l

連絡先:本武陽一,東京大学大学院総合文化研究科, mototake@sacral.c.u-tokyo.ac.jp

次元の第 *l* 層の出力とし, *f*(·) を,

$$f(w^{l-1}h^{l-1} + b^{l-1}) = (f_1, f_2, \cdots, f_{d_l})$$

$$f_j = f\left(\sum_{i}^{d_{l-1}} \left(w_{ij}^{l-1}h_i^{l-1} + b_j^{l-1}\right)\right)$$
(2)

とした. *f* は活性化関数と呼ばれ,良くシグモイド関数や LeRU 関数等が用いられる.従って,ある入力 *p* が与えられた場合 の DNN 写像関数のヤコビアン行列 *J* は,

$$J_{ij}^{\text{DNN}} = \frac{\partial F_j(\boldsymbol{x})}{\partial x_i}$$
$$= \sum_{\langle k_{l-1}, k_{l-2}, \cdots, k_1 \rangle} \frac{\partial h_j^L}{\partial h_{k_{l-1}}^{l-1}} \frac{\partial h_{k_{l-1}}^{l-1}}{\partial h_{k_{l-2}}^{l-2}} \cdots \frac{\partial h_{k_1}^1}{\partial x_i}$$
(3)

で与えられる.また、各層間の写像関数のヤコビアンを

$$\boldsymbol{J}^{l} = \begin{pmatrix} \frac{\partial h_{1}^{l}}{\partial h_{1}^{l-1}} & \cdots & \frac{\partial h_{1}^{l}}{\partial h_{d_{l-1}}^{l-1}} \\ \vdots & \ddots & \vdots \\ \frac{\partial h_{d_{l}}^{l}}{\partial h_{1}^{l-1}} & \cdots & \frac{\partial h_{d_{l-1}}^{l}}{\partial h_{d_{l-1}}^{l-1}} \end{pmatrix}$$
(4)

とした場合に,全体のヤコビアンは

$$\boldsymbol{J}^{\text{DNN}} = \boldsymbol{J}^{L} \boldsymbol{J}^{L-1} \cdots \boldsymbol{J}^{1}$$
(5)

で与えられる. 例えば活性化関数 *f* としてシグモイド関数を 用いた場合,

$$\boldsymbol{J}^{l} = \begin{pmatrix} h_{1}^{l-1}(1-h_{1}^{l}) & \dots & h_{d_{l-1}}^{l-1}(1-h_{1}^{l}) \\ \vdots & \ddots & \vdots \\ h_{1}^{l-1}(1-h_{d_{l}}^{l}) & \dots & h_{d_{l-1}}^{l-1}(1-h_{d_{l}}^{l}) \end{pmatrix}$$
(6)

となる. このように写像関数から定義された多様体上のある点 p での局所近傍座標系である接空間 T_pM は、点 p 周りのでヤ コビアン行列 J の縮退していない基底空間であるユークリッ ド空間として与えられる.

この描像は,DNN がシグモイド関数のほぼ線形な部分をつ なぎ合わせた曲がった空間上で主成分分析のようなことを行う



図 1: DNN におけるデータ分布の多様体への引き込みメカニ ズムの模式図

ものとして理解される.その際圧縮される方向は,多様体から 遠く離れた点まで及ぶことに注意されたい(図2).また,シ グモイド関数の線形部分が有限であることから,多様体から離 れるほど,圧縮領域の隙間となる空間が拡大する.このような 空間は,シグモイド関数の性質から,0あるいは1に吸引され ていくと考えられる.

2.2 DNN のモデル化した多様体からの対称性抽出手法

前節の議論より,入力空間で多様体上にないデータ点は,多 様体上に吸引されるとわかる.入力空間から見た中間層での データ点の移動は図2のような模式図に従うと考えられる.中 間層において一度多様体上に吸引されたデータ点は,それより も後の層において多様体上のデータと分離されることはない. したがって,出力層 **F**(**x**)において,全てのデータは多様体上, あるいは空間の角となる0と1の組み合わせで与えられる座 標点に縮退する.

これをふまえると、あるデータ多様体を学習した DNN に対して、新たに与えられたデータセット $\{x_i\}_{i=1}^{N}$ がそのデータ 多様体上にあるかは、それが出力空間に写像された際のデータ 集合 $\{F(x_i)\}_{i=1}^{N}$ との間の二乗誤差

$$E(\{\boldsymbol{x}_i\}_{i=1}^N; \boldsymbol{F}(\boldsymbol{x})) = \sum_{i=1}^N [\boldsymbol{x}_i - \boldsymbol{F}(\boldsymbol{x}_i)]^2$$
(7)

によって評価されるとわかる (図 2).

なぜならば,入力空間で多様体から離れたデータほど,出力 空間で多様体上に落ち込むまでに移動した距離が遠く,結果と して二乗誤差が大きくなるためである.

今知りたいのは、DNNの訓練時に用いたデータ多様体分布が どのような座標変換 A に対して普遍となるかということである. 先ほど議論した多様体からの距離を測る指標 $E({x_i}_{i=1}^N; F(x))$ を用いると、データ多様体の分布構造を変えない座標変換は、

$$E(\{\mathbf{A}\boldsymbol{x}_i\}_{i=1}^N; \boldsymbol{F}(\mathbf{A}\boldsymbol{x})) = \sum_{i=1}^N \left[\mathbf{A}\boldsymbol{x}_i - \boldsymbol{F}(\mathbf{A}\boldsymbol{x}_i)\right]^2 \qquad (8)$$

が小さな値となる A を集めてくれば良いとわかる.

具体的には、行列 **A** の $d_{in} \times d_{in}$ 個ある要素 a_{jk} について サンプリングを行うことで、 $E(\{\mathbf{A}\boldsymbol{x}_i\}_{i=1}^N; F(\mathbf{A}\boldsymbol{x}))$ がある閾値 以下となる **A** の集合が得られる.

2.3 ネーターの定理

系の対称性と保存量との関係について、ネーターによって証明された定理である [9]. ここで、d 次元一般化座標を、 $\mathbf{q}(t)$ と $\mathbf{p}(t)$ とし、系のハミルトニアンを H(q, p; t) とする、今、ハ ミルトニアン H(q, p; t) が、ある微小な座標変換

$$t' = t + \delta t$$

$$q'_i = q_i + \delta q_i$$
(9)

(ただし, j=1 ~ d) に対して普遍であるとする. すると, ネー ターの定理より微小変換の生成子 G_{δ} が時間によらず普遍な量 となる. つまり, 保存量となる. ここで微小変換の生成子 G_{δ} は,

$$(\delta q_j, \delta p_j) = \left(\frac{\partial G_\delta}{\partial q_j}, \frac{\partial G_\delta}{\partial p_j}\right) \tag{10}$$

を満たすものとして定義される.



図 2: 対称性の抽出. 多様体上からの距離を,入力データと出力データの2 乗誤差から定量化.

2.4 時系列データの対称性とネーターの定理

次に、ネーターの定理を適用するために必要なハミルトニ アン H(q, p; t) の微小変換に対する普遍性を, 位相空間での時 系列データ $\{\mathbf{q}^{i}(t), \mathbf{p}^{i}\}_{i=1}^{N}$ から得る方法について述べる. ハミ ルトニアン力学系において、力学系の時間発展方程式は、q,p についての一階微分方程式となる. このことから, 一つの仮 定を設定した. それは、「与えられたハミルトニアン H(q, p; t) の元である一定エネルギー E で運動を行う系が、4d 次元空間 (q(t), p(t), q(t+1), p(t+1))に取り得る全ての状態が持つ等 エネルギー面(多様体 M_E)の持つ対称性が,ハミルトニア ンの対称性と関係する.」というものである.これはすなわち, あらゆる初期値から生成した同じエネルギーを持つ運動の時系 列データの持つ多様体の対称性を抽出すれば,保存量の候補が 得られるという仮定となっている. この仮定によって得られた 結果の妥当性は、最終的に得られた保存量が時間普遍であるこ とを確認することで容易に確認される. この M_E が持つ対称 性を 2.2 節の方法で推定することで、保存量推定を行うのが本 提案手法の枠組みである.

2.5 変換の生成子

最後に, 2.2 節の方法で得られるハミルトニアン普遍な大域 的変換 A から, 微小変換の生成子 G_{δ} を取り出す方法を, 3 次 元回転の行列と SO(3) の生成子を例にとって説明する.

Z 軸まわりの回転行列を書き下すと,

$$A_z = \begin{pmatrix} \cos(\theta_z) & \sin(\theta_z) & 0\\ -\sin(\theta_z) & \cos(\theta_z) & 0\\ 0 & 0 & 1 \end{pmatrix}$$
(11)

となる. 今, この回転角 θ_z が, 非常に小さい (= $\delta \theta_z$) とし て, この行列を θ_z についてテーラー展開し, 一次の項までを とると,

$$A_z \sim \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} + \begin{pmatrix} 0 & -\delta\theta_z & 0 \\ \delta\theta_z & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$
(12)

となる.この右辺第二項が Z 軸周りでの回転の生成子となる. 他の軸についても同様に生成子が得られる.

実際にサンプリングして得られる A_z においては, 隠れ変数である θ_z による変換行列のパラメトライズはできないため.

要素 a_{ij} の間の関係の回帰や,それらの分布の分析などから,

周りでの微小変動の方向を見積もることで、回転の生成子を 得た.

結果と考察

物理学データへの手法の適用を行う.具体的に適用するの は、以下のハミルトニアン

$$H(\mathbf{x}, \mathbf{p}) = \frac{\mathbf{p}^2}{2m} - \frac{k}{r} \tag{14}$$

に従う 2 次元空間中の中心力運動である.実験では, k = 10,m = 1 とし,さらに,運動を円運動に限定した上で,系の 取りうる全ての状態によってトレーニングデータを生成した. サンプリング対象となる変換行列は,配位空間 x = (x1,x2)

ž,

$$\begin{pmatrix} x'_1 \\ x'_2 \end{pmatrix} = \mathbf{A} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} a_{11} & a_{21} \\ a_{12} & a_{22} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \quad (15)$$

と変換する.この時,運動量空間 $\mathbf{p} = (p_1, p_2)$ は,

$$\begin{pmatrix} p_1' \\ p_2' \end{pmatrix} = \begin{pmatrix} m\dot{x'}_1 \\ m\dot{x'}_2 \end{pmatrix} = \begin{pmatrix} a_{11} & a_{21} \\ a_{12} & a_{22} \end{pmatrix} \begin{pmatrix} m\dot{x}_1 \\ m\dot{x}_2 \end{pmatrix}$$

$$= \begin{pmatrix} a_{11} & a_{21} \\ a_{12} & a_{22} \end{pmatrix} \begin{pmatrix} p_1 \\ p_2 \end{pmatrix}$$
(16)

と変換される.したがって、この4つの変換行列の要素につい てサンプリングを行う.ちなみに今回は、時間変換に対する対 称性は、明らかであるので検討対象としなかった.ちなみに、 時間方向の対称性は、エネルギー保存則を導く.

ここではメトロポリスサンプリングを用いた.メトロポリス サンプリングでは,全てのパラメータについて step 数は 0.2 と設定され,分散 0.2 のガウスノイズを仮定した.その上で, 1,000step の burn in の後に 9,000step のサンプリングを実行 した.この結果,図 3 のような結果が得られた.



図 3: 中心力ポテンシャル中を運動する物体の時系列データ セットに対して,提案手法によって得られたデータ多様体の分 布を変えない対称性変換 A の要素 *a_{jk}* の比較. A が回転行列 となっている場合に得られるであろう分布構造を黒線で表した. ただし, *a*₁₁-*a*₂₂ と *a*₁₂-*a*₂₁ の図における直線の片方は, 鏡像 変換の従う分布構造についてもプロットした図となっている.

得られた結果より、変換の生成子を推定する.今回得られた結果の図から変換が単位行列となる $a_{11} = a_{22} = 1$ 周りの 微小変動を考察すると、 $a_{11}-a_{12}$ と $a_{11}-a_{21}$ の関係(図3)より、生成子の a_{11},a_{22} 成分が0となり、 a_{12} と a_{21} 成分が $\pm \delta \theta_y$ 方向を持つとわかる.さらに、 $a_{21}-a_{12}$ の関係より、生成子の $a_{11,a_{22}}$ 成分の符号が逆であることがわかる.以上より、得られた変換の生成子は、

$$\begin{pmatrix}
0 & -\delta\theta_y \\
\delta\theta_y & 0
\end{pmatrix}$$
(17)

となる. これはちょうど SO(2) の対称性の生成子となっており、これより、角運動量が保存することがわかる.

以上のように,提案手法によって力学系の時系列データから,そのハミルトニアンの対称性に関する情報を抽出し,保存 則が推定されることが確認された.

4. まとめと今後の展望

データ多様体をモデル化した DNN から,その多様体の持 つ対称性に関する情報を抽出する手法を構築した.そしてその 応用として,中心力ポテンシャル中を運動する力学系の時系列 データの作るデータ多様体に提案手法を適用した.その結果, 力学系が持つ対称性の抽出と,ネーターの定理を介した保存量 の推定を実現した.

本提案手法は、さらに複雑な対称性を持つ系への適用も可 能であると考えられる.具体的には、中心力ポテンシャル中の 力学系が持つ配位空間とは違う空間における隠れた対称性であ る、SO(4)の抽出と、対応する保存量であるルンゲ-レンツベ クトルの推定等の実現が期待される.

近年の DNN 技術の有用性の確認とともに,多数の研究者・ 技術者が各種の DNN アルゴリズムの開発やパラメータチュー ニングを行なっている.この状況は,各種データセットに対す る多様体構造についての莫大な知見が蓄積されつつあることを 意味する.本提案手法は,このような蓄積された莫大な知見を 解析するための手法開発に繋がっていくことが期待される.

参考文献

- Irie Bunpei and Kawato Mitsuo. Acquisition of internal representation by multi-layered perceptrons. The Transactions of the Institute of Electronics, Information and Communication Engineers D, Vol. 73, No. 8, pp. 1173–1178, 1990.
- [2] GE Hinton and RR Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 2006.
- [3] Pratik Prabhanjan Brahma, Dapeng Wu, and Yiyuan She. Why deep learning works: A manifold disentanglement perspective. *IEEE transactions on neural networks and learning systems*, Vol. 27, No. 10, pp. 1997– 2008, 2016.
- [4] Ronen Basri and David Jacobs. Efficient representation of low-dimensional manifolds using deep networks. arXiv preprint arXiv:1602.04723, 2016.
- [5] Salah Rifai, Yann N Dauphin, Pascal Vincent, Yoshua Bengio, and Xavier Muller. The manifold tangent classifier. pp. 2294–2302, 2011.
- [6] 本武陽一,池上高志. Deep neural networksの力学的・幾 何学的解析. 人工知能学会全国大会論文集, Vol. JSAI2016, pp. 1A5OS27c1-1A5OS27c1, 2016.
- [7] Yhoichi Mototake and Takashi Ikegami. The dynamics of deep neural networks. Proceedings of the Twentieth International Symposium on Artificial Life and Robotics, Vol. 20, , 2015.
- [8] Y. Mototake. Geometrical structures embedded in high dimensional data sets and deep learning : Analysis and application to dynamical systems. 2016.
- [9] AE Noether. Nachr kgl ges wiss göttingen. Math. Phys. KI II, Vol. 235, 1918.