畳み込みニューラルネットワークと全方位空間画像による都市空間の特徴抽出と評価手法

○衣川 雛*1, 瀧澤 重志*2

キーワード:都市景観,全方位画像,深層学習,CNNオートエンコーダー,CAM,次元圧縮

1. はじめに

空間の印象を評価するために,建築や都市計画の分野 では,画像を用いた分析が広く行われている.例えば, 空の見える量を示す天空率 ¹⁾や緑の量を示す緑視率 ²⁾な どは比較的よく用いられている.一般に画像を使って空 間の評価を行う場合は,画像から何かしら意味のある計 算可能な指標を抽出する必要があると考えられている. しかし,空間の印象を左右する特徴は無数に考えられ, 明示的に定義できるとは限らない.

近年,画像から特徴量を自動的に抽出する深層学習を 使った、景観等の印象評価の研究が行われている、それ らには例えば, 矢吹ら³⁾, 阿部ら⁴⁾, 山田ら⁵⁾, Stephen⁶⁾ ら, Liu ら ⁷⁾のものなどがある. 筆者らも既往研究 ⁸⁾で, 深度情報を含む3次元の全方位画像を、ゲームエンジン を用いて仮想空間からリアルタイムに取得する空間評価 システムを開発した.そして、VRを用いた印象評価実験 を行い、その評価結果を CNN ベースの画像判別モデル で学習し、検証データの予測精度を把握し、提案手法の 有効性を検証した.この研究では、画像から特徴量を事 前定義することなく学習で抽出することになるが、目的 変数を設定した上でモデル全体を学習させる、すなわち 教師あり学習を行うので、目的変数が変化したら、ファ インチューニングなどの方法はあるにせよ、原則、都度 モデル全体を学習させる必要がある.一般に深層学習の 学習には時間がかかるし、場合によってはニューラルネ ットワーク以外の学習モデルを使って、判別規則を明示 化するといった用途も考えられるので、現状の教師有り 学習の枠組みでは、画像データを都市の多様な分析に生 かすことには限界があると考えられる.

そこで本研究では,教師無し学習の方法である CNN オ ートエンコーダーを使って,事前に全方位画像の高次元 情報を目的変数に依らず次元圧縮し,その後,以前行っ た印象評価実験の結果を学習させ,既報のモデルと精度 や注目領域部分の差異を比較する.

2. 既報⁸⁾の内容のまとめ

本研究では3次元空間評価システムの基盤として、ゲ ームエンジンとして普及が進んでいる Unity を用いる. 以下に本研究に関連する既報の内容をまとめる.

2.1 対象空間の設定

ゼンリンから公開されている Unity 用の3次元空間デ

ータ ZENRIN City Asset Series⁹⁾の難波エリア (Japanese Naniwa City) を評価対象の空間とした. このエリアは難 波エリアを中心とした,東西約 700m 南北約 600m の範 囲をモデル化したものである (図 1).



図 1 Unity 用の 3 次元マップ: Japanese Naniwa City⁹⁾



図2 観測点(黒〇)と評価点(緑〇)

印象評価実験では、歩行者が普段歩く場所から街を眺めた際の印象を質問するので、道路上の歩行者が歩く場 所で空間特徴量を計測する必要がある.場所の選定の基 準として、歩道がある区域では歩道幅の中心、歩道がな く車道で交通量が多い地域では,道路の両端、歩道がな く交通量が少ない地域では道路の中央を歩行者の通行場 所と仮定し、それらの場所を5m間隔で座標を収集した. その結果、図2に示す2,029点の観測点が得られたが、 それらの点をすべて使って印象評価実験を行うのは困難 なので、それらの中から図2に示す50地点を無作為に 選択し、印象評価実験における評価点とした.

2.2 空間特徴量の抽出

各観測点の特徴を抽出するために、Unity のアセット である Spherical Image Cam (現在は販売終了)を利用し、 カメラが存在する場所の全方位画像を、リアルタイムで 撮影・保存できるように、Unity上でスクリプトを作成し た.全方位の RGB 画像の例を図 3 に示す.なお、既報で は RGB 画像の他に深度画像も出力していたが、本研究 では(時間の制限により) RGB 画像のみで学習を行った.



図3 全方位画像の例

2.3 VRによる印象評価実験システムの構築と実験

Unity に Oculus Rift を接続し、印象評価実験を行うためのインターフェースを作成した(図4).



図4 VR による印象評価実験の様子

2016年の冬に建築系の大学生・大学院生8名を対象として印象評価実験を行った.評価は1(悪い),2(少し悪い),3(少し良い),4(良い)の4段階を,直感で判断するよう教示を行った.いずれの被験者も,評価値が4の場所は少なく,2前後が多かった.そのため4クラス分類は難しいので,評価値が1,2の場所をClass0,評価値が3,4の場所をClass1とした2分類問題として学習を行うことにした.表1に2クラス分類の場合の各クラスの度数分布を示す.全体的に低評価の地点の方が多い.

表1 被験者ごとの評価ク	ラスの分布
--------------	-------

評価値	被験者ごとのクラスの集計							
	А	В	С	Ð	Е	F	G	Н
1 or 2: Class 0	42	34	23	29	30	40	32	26
3 or 4: Class 1	8	16	27	21	20	10	18	24

2.4 CNN による学習

深層学習のライブラリとして Chainer¹⁰⁾を用い, そのサ ンプルファイルの一つである train_imagenet.py をカスタ マイズし, GoogLeNet¹¹⁾ (Batch normalization 版) を用い て学習を行った.



図 5 GoogLeNet の構造¹¹⁾

Unityを用いて、学習に使う画像を各地点で撮影する. 画像を扱う深層学習では、人工的にできるだけ多くのデ ータを用意する必要がある.本研究では各観測点につい て、図6に示すように、カメラを鉛直軸に対して0度か ら20度刻みで340度まで回転させて、同一地点で18枚 ずつ撮影した.通常の全方位画像は縦横比が縦:横=1: 2の正距円筒図法だが、CNNに入力するデータの制限か ら、画像サイズを256×256ピクセルの正方形としてJPG 形式で出力した.



図6カメラの回転による同一地点の全方位画像の増量

さらに、学習用と検証用データがおおよそ3:1になるように、各クラスの比率を考慮して、被験者ごとにランダムに地点を分割した.以上の設定に基づき、各被験者で学習用の画像データセットでDCNNの学習を行い、得られたモデルに対して検証用の画像データセットを適用して、それらの誤答率を調べた.

3. CNN オートエンコーダーによる画像情報の次元圧縮

ここからは、本研究で行った CNN オートエンコーダ ー(以下 CNNAE と訳す)について説明する.基本とな るオートエンコーダーは、高次元の入力データを何層か の隠れ層で次元圧縮した後(Encode)、そこから反対に高 次元化(Decode)して、元のデータとできるだけ同じよ うな出力を得ようとする手法である.これを画像処理デ ータでできるように、畳み込み層で行うようにしたのが CNNAE である.図7に CNNAE の概念図を示す.入力 画像に関して畳み込み演算(CONV)が行われ、中間の最 も画像サイズが小さくなる層が特徴表現層である.そこ から逆畳み込み(DCONV)を行って、元の画像にできる だけ類似した画像を出力できるように、各層のパラメー タを学習させる.特徴表現層では、大きなサイズの入力 画像が、多数の小さな画像に変換されている.この層の 情報を使って、画像の特徴ベクトルを抽出できる.具体 的には、特徴表現層の各チャンネルの画像全体を平均プ ーリング(GAP)することで、フィルタ数の次元の特徴 ベクトルを得ることができる.このようにして得た特徴 ベクトルを、本研究では単純パーセプトロンによって、 次に述べる評価実験の学習に用いる.



4. 評価実験結果の学習

表2に学習で用いる CNNAE の緒言を示す.表2に示 すように4つのモデルを作り,2.1 で説明した2,039 個の 観測点で撮影された全方位画像16,232 枚について,それ ぞれの CNNAE を学習させた.

表 2 実験で用いる CNNAN の緒言						
	CNNAE1	CNNAE2	CNNAE3	CNNAE4		
層数	3	3	2	2		
フィルタ数	16, 32, 40	8, 16, 32	16, 32	16, 32		
フィルタ幅	5, 5, 3	5, 5, 3	5, 5	8,8		
ストライド	3, 3, 2	3, 3, 2	3, 3	4,4		
特徴表現層の	12	12	27	14		
画像サイズ	13	12	27	14		

衣 う アストアータの正解率					
被験者	GoogLeNet	CNNAE1	CNNAE2	CNNAE3	CNNAE4
А	0.71	0.85	0.85	0.85	0.85
В	0.80	0.69	0.71	0.77	0.72
С	0.55	0.69	0.48	0.62	0.58
D	0.81	0.75	0.65	0.72	0.72
E	0.64	0.68	0.72	0.59	0.70
F	0.96	0.85	0.87	0.90	0.85
G	0.86	0.75	0.78	0.73	0.68
Н	0.77	0.75	0.70	0.61	0.61
Ave.	0.76	0.75	0.72	0.72	0.71
Std.	0.12	0.06	0.11	0.11	0.09

次に,50の評価点の画像900枚について,学習したそ れぞれの CNNAE を使って特徴量を抽出し,2章の分類 問題のデータを,単純パーセプトロンを使って学習させ, 10-fold 交差検証により,平均正答率を求めた.それらの 結果と,2.で行った過去の GoogLeNet の正解率を比較し たものを表3に示す.最も構造が複雑な CNNAE1 の平均 精度は GoogLeNet に近く,その標準偏差は GoogLeNet を 凌駕している.精度については,CNNAE による教師無し 学習でも,教師あり学習のモデルを使った場合と比較し て安定して良好な結果を得ることができた.

5. Class Activation Mapping(CAM)¹²⁾による注目領域の 比較

CAM は、CNN ベースのニューラルネットで分類問題 を解いた時に、分類モデルが画像のどこに注目したのか を可視化する手法である.この概念図を図 8 に示す. CAM では CNN の最終層の画像の画素値を平均プーリン グして、得られた特徴ベクトルの重み付き和として分類 モデルを構築する.最終層の特徴マップの画素をその重 みで積を取った上で、和を取って拡大することで、CNN の注目地点を可視化できる.



図 8 CAM の概念図¹²⁾

この方法を使って,被験者 Fの GoogLeNet と CNNAE1 モデルの注目領域を可視化した例を図 9,10 に示す. GoogLeNet の場合,高評価地点では,街路樹のあたりに はっきりとした注目領域が形成されている.一方,低評 価地点の注目領域は奥行の方向に延びている. CNNAE1 の場合,高評価地点では木や手前の建物の1階のファサ ードなどに,低評価地点では全体的に注目領域が広がっ ている. CNNAE の方は注目領域が分散する傾向があり, GoogLeNet と比較して解釈が難しいと結論付けられる.

6. おわりに

本研究では,教師無し学習の方法である CNNAE を使 って,事前に全方位画像の高次元情報を目的変数に依ら ず次元圧縮し,その後,以前行った印象評価実験の結果 を学習させ,既報のモデルと精度や注目領域部分の差異 を比較した.その結果,精度や安定性は教師有学習と比 較して遜色ないことを確認した.一方,CNNAE の注目領 域は解釈が難しいため,その理由はモデルの構造などさ らなる分析が必要と考えられる.



(a) 元画像

(b) GoogLeNet 図 9 CAM による CNN の注目点(高評価地点)

(c) CAE1



(a) 元画像

(b) GoogLeNet 図 10 CAM による CNN の注目点(低評価地点)

(c) CAE1

謝辞

本研究の一部は科研費基盤 C(16K06652), 基盤 A(V16H01707) の補助の下で行われました.研究を主に担当した,当時学部 4年生の古田愛理さんと上田健之祐君に感謝します.

[参考文献]

- 国土交通省住宅局:建築物に対する景観規制の効果の分析 手法について,2007, http://www.mlit.go.jp/jutakukentiku/jutakukentiku.files/keikan/keikankisei.pdf
- 2) 藤井健史,山田悟史,廣瀬徳郎,及川清昭:CG モデルによ る全方位緑視率の計量手法とその応用可能性,日本建築学 会技術報告集,19(43), pp.1067-1072, 2013.
- 3) 矢吹和也, 安福健祐, 阿部浩和: Deep Learning を用いた景観 評価の手法に関する基礎的研究, 日本図学会大会学術講演 論文集, pp.23-28, 2015.
- 4) 阿部浩和, 李ロウン, 安福健佑: 街路空間評価におけるディープラーニングの適用可能性, 日本図学会大会学術講演 論文集, pp.85-90, 2016.
- 5) 高橋秀彬,山田悟史: Deep Learning を用いた街並み画像の 分類と感性評価の推定,日本建築学会第40回情報・システ ム・利用・技術シンポジウム論文集:報告, pp.329-329, 2017.
- S. Law et al.: An application of convolutional neural network in street image classification: the case study of London, ACM GeoAl'17 Proc. 1st Work. on Artif. Intell. Deep. Learn. for Geogr.

Knowl. Discov, 2017.

- L. Liu, et al.: A machine learning-based method for the large-scale evaluation of the qualities of the urban environment, Computers, Environment and Urban Systems, 65, pp.113-125, 2017.
- 8) A. Takizawa and A. Furuta, 3D Spatial Analysis Method with First-Person Viewpoint by Deep Convolutional Neural Network with Omnidirectional RGB and Depth Images, eCAADe 2017, Sapienza University of Rome, Rome, Italy, pp.693-702, 22 Oct. 2017.
- 9) Zenrin 株式会社: ZENRIN City Asset Series, <u>http://www.zenrin.co.jp/product/service/3d/asset/</u>
- 10) S. Tokui et al.: Chainer: a Next-Generation Open Source Framework for Deep Learning, Proceedings of Workshop on Machine Learning Systems in The Twenty-ninth Annual Conference on Neural Information Processing Systems, 2015.
- 11) A. Krizhevsky et al.: ImageNet Classification with Deep Convolutional Neural Networks, Pereira, F. et al. (eds), Advances in Neural Information Processing Systems 25, Curran Associates, Inc.,pp.1097-1105, 2012.
- 12) B. Zhou, et al.: Learning Deep Features for Discriminative Localization, Computer Vision and Pattern Recognition, 2016.

*1 大阪市立大学生活科学部 学部生

*2 大阪市立大学生活科学研究科 教授 博士 (工学)