

An unsupervised approach to geographical knowledge discovery using street level and street network images

Stephen Law*
The Alan Turing Institute &
University College London
London, UK

Mateo Neira*
The Alan Turing Institute &
University College London
London, UK

ABSTRACT

Recent researches have shown the increasing use of machine learning methods in geography and urban analytics, primarily to extract features and patterns from spatial and temporal data using a supervised approach. Researches integrating geographical processes in machine learning models and the use of unsupervised approaches on geographical data for knowledge discovery had been sparse. This research contributes to the ladder, where we show how latent variables learned from unsupervised learning methods on urban images can be used for geographic knowledge discovery. In particular, we propose a simple approach called Convolutional-PCA (*ConvPCA*) which are applied on both street level and street network images to find a set of uncorrelated and ordered visual latent components. The approach allows for meaningful explanations using a combination of geographical and generative visualisations to explore the latent space, and to show how the learned representation can be used to predict urban characteristics such as street quality and street network attributes. The research also finds that the visual components from the *ConvPCA* model achieves similar accuracy when compared to less interpretable dimension reduction techniques.

CCS CONCEPTS

• **Computing methodologies** → **Unsupervised learning**; *Computer vision*; • **Applied computing** → **Architecture (buildings)**.

KEYWORDS

urban analytics, unsupervised learning, convolutional neural networks, knowledge discovery, computer vision, machine learning

1 INTRODUCTION

According to [22], Geographic knowledge discovery (GKD) is the process of using computational methods and visualisation to explore spatial databases to discover useful geographic knowledge. Despite the ubiquity of geographically-labelled image data and the subsequent use of machine learning methods to retrieve geographical information, the majority of the researches have focused mainly on the use of supervised learning approaches. For example, on the use of convolutional neural networks *CNN* to make inferences on perceived safety [24], house price [19] and scenicness [29]. These researches required effort on both collecting the data and on learning a specific objective. As a result, there is an opportunity to use urban image information in an unsupervised and scalable manner. Our research question therefore implies, what compact latent representation can be learnt from urban images without supervision, how

can this information be described and what is this representation useful for?

This study contributes to these research questions and proposes an unsupervised learning model called Convolutional-PCA (*ConvPCA*) that summarises urban imagery into a set of lower dimensional uncorrelated latent components. We apply this method to two case studies namely for: Google StreetView images [10] and OpenStreetsMaps (*OSM*) street network images [26]. In the experiments, we first map and visualise the extremes of the responses geographically and generate new synthetic images by perturbing the values of each component whilst holding all the other component values constant. We then study the latent components by using it to predict different geographical datasets such as street enclosure and street frontage type for the StreetView image data and network density and network centrality for the street network image data. The research finds that the visual components from the *ConvPCA* model have interpretable meanings with predictive abilities to geographical labelled data using a compact representation. The research also finds that the visual components from the *ConvPCA* model achieve a similar accuracy to other dimension reduction techniques such as an autoencoder while retaining its interpretability. From a machine learning perspective, we gain new knowledge about these latent components which contribute to the recent efforts in linking the two disciplines [16] [27].

2 RELATED WORKS

2.1 StreetViews

Street-level images have been used extensively in intelligent transportation systems research. Specifically on the deployment of autonomous vehicles where deep convolutional neural networks (*CNN*) had been applied for urban scene understanding [28]. More recently, we have also seen the use of generative models such as Generative Adversarial Networks (*GAN*) to synthetically create street scenes that can be used to train self-driving vehicles [33]. Despite its popularity in transportation research, there had been limited effort on using street-level imagery to retrieve geographical information and for studying urban planning problems. One such example is StreetScore where [24] collected subjective human perception data from street images through a crowd-sourced survey (Place Pulse 2.0) which are then used to predict the perceived safety of a place [8]. Another example is the work of Gebru et al. [9] whom extracted features such as car types from Google StreetView images to predict the income, race, education, and voting patterns for cities in the United States. We have also seen the use of urban images [29] to predict scenicness ratings which were found to affect urban wellbeing as well as the type of urban frontages [19] which

*Both authors contributed equally to this research.

is an important urban design attribute. These fairly recent efforts relied on extracting visual features from street-level images which are then related to different socio-economic factors. In contrast to these works, Law et al [18] extracted a visual response from urban images by directly estimating house prices. A distinguishing difference here is that the method extracted a visual scalar response that corresponds directly to house price, which can be visualised and interpreted in traditional econometric model. In summary, these recent researches focused on learning a set of visual features or response from urban imagery using a supervised learning approach. Our research extends from this work where we propose a two stage method in learning a set of generic and compact latent visual components from an unsupervised learning approach. We then, through a set of analysis and experiments interrogate, describe and explore these components for geographic knowledge discovery.

2.2 Street Networks

In the case of street networks, there has been a long-standing effort to analyse and to understand them from a quantitative perspective and to generate models that are able to reproduce their empirical features. Previous works have largely been based on complexity theory and network science [5, 20, 31]. This includes analyzing the spatial configuration of urban street networks [13] and analyzing urban systems from an information theoretic perspective [2].

More recently, there has been a growing interest in applying machine learning methods to extract useful information from the vast amount of data now openly available from sources such as *OSM*. Examples of such works have used neural networks to classify street network patterns of different cities, where two different methods had been used. The first used a Convolutional Autoencoder *CAE* to create dense urban vectors that are used to cluster similar urban morphologies using a self-organizing map [23]. The second approach used a Variational Auto Encoder *VAE* to measure similarity across different networks [14].

Generative models have also been used to generate synthetic street networks. Variational Autoencoder trained on street network images has been used by sampling from the latent space z [14], however the resolution of these are low, and fail to capture fine grain detail of local streets. A Generative Adversarial Network such as *StreetGAN* [12] has also been proposed to generate a multitude of arbitrary sized street networks that faithfully reproduce the style of the original dataset.

Current limitations in the use of VAE, CAE, and GANs on street networks lie in the interpretability of the latent space and its relationship to geometrical and topological properties used in established network measures. Our research contributes to these researches by developing a methodology to interpret the lower-dimensional embedding learnt by a convolution autoencoder. This allows for greater interpretation of the unsupervised model, as well as providing some initial results as to the relationship between the embedding and the established network measures.

3 METHODS AND MATERIALS

3.1 Convolutional-PCA

We propose here the Convolutional-PCA (*ConvPCA*), which combines a type of Convolutional Neural Network called the Convolutional Autoencoder (*CAE*) with a linear PCA (PCA_{lin}) to retrieve a set of latent visual components that summarise a StreetView image or a street networks image. We first describe the *CAE* followed by the PCA_{lin} . Deep Convolutional Autoencoder *CAE* is an unsupervised method that uses convolutional neural network (*CNN*) to learn a compact representation or a set of visual features [3, 11, 21]. Deep *CAE* consists of two set of layers, an encoder $f_w(\cdot)$ and a decoder $g_u(\cdot)$

$$f_w(x) = \sigma(x * W) \equiv z \quad (1)$$

$$g_u(z) = \sigma(z * U)$$

where x is the input vector, z is the latent features, $*$ is the convolution operator that extract image features and σ is a *ReLU* activation function to model nonlinearity in the neural network. These convolutional layers can be stacked sequentially where the encoding layers reduce the dimension to a latent variable z while the decoding layer increases the dimension back to image space. The sequential architecture can be seen in figure 1. Following [21], the parameters of the encoder $z = F_w(x)$ and the decoder $x' = G_u$ are updated by minimising the reconstruction losses between x and $x' = G_u(F_w(x_i))$.

$$L_r = \frac{1}{n} \sum (x_i - G_u(F_w(x_i)))^2 \quad (2)$$

In our research, we further compress the latent visual features by applying a linear principal component analysis PCA_{lin} which summarises the visual feature z into a set of linearly uncorrelated variables v . To compute v , we first standardise z and compute the eigenvectors and eigenvalues of the feature covariance matrix P . We then take the eigenvectors to calculate the full principal component decomposition of z , given by $V = XW$, where W is the eigenvector matrix. V can be re-projected back on to the original latent space produced by the encoder before passing in to the decoder to reconstruct the images. This process allows us to:

- Retrieve a set of uncorrelated and ordered visual latent components that can be visualised and mapped geographically.
- Make changes to individual components and decoding it to generate a synthetic image.
- Relate learnt visual latent components to geographical labelled data.

To discover new geographical knowledge and in testing the usefulness of the latent representation, we will visualise these components by generating new images when perturbing in the PCA_{lin} space and also in mapping them. We will then use these components for down stream tasks such as prediction and classification. The process can be seen in figure 1. Further research is required to validate the meaning of these visual latent components quantitatively and in comparing it with latent components extracted from other methods. These limitations will be further elaborated in the discussion section.

PCA_{lin} is selected as it is a well principled dimension reduction technique that learns a compact and meaningful representation with uncorrelated and ordered components. Approaches such as autoencoders can find a similar representation but without the same interpretability [17]. In the prediction experiment, we will study and compare the extent a PCA_{lin} , a linear autoencoder and a non-linear autoencoder are able to learn a compact representation for different down-stream tasks. A benefit of finding uncorrelated ordered components is that these factors can be inserted into a generalised linear modelling framework, whose coefficients can be interpreted. Such research are not explored in this paper but the coefficients are meaningful for example in econometric studies [18].

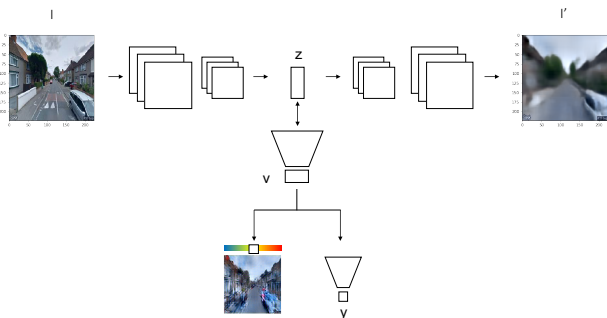


Figure 1: Architecture of ConvPCA, which combines a Convolutional AutoEncoder (CAE) with a linear PCA (PCA_{lin}) to retrieve a set of uncorrelated and ordered visual latent components that summarises street-level and street-network image data. These latent components are then used for knowledge discovery through visualisation and different down-stream tasks such as prediction and classification.

3.2 Materials

We collected two datasets. The first dataset is street images taken from the Google StreetView API [10]¹. Similar to [19], we collected a front-facing image for each street in the Greater London Area. To collect the dataset, we constructed a line-graph from the street network of London (OS Meridian line2 dataset [32]). We then take the geographic median and the azimuth of the street edge to give both the location and the bearing when collecting each image. We collected a total of 110,493 street images in London. For more details in the data collection method please see [19]. Figure 3 illustrates typical images from the dataset.

The second dataset is the street network dataset taken from OpenStreetsMaps [26], we query all the cities and towns for a total of 107,973. For each city and town we download the street network within a 1.5km x 1.5km box at the centroid of each place using osmnx [4], as shown in Figure 4. For each 1.5km x 1.5km grid we retrieve a graph $G = (V, E)$ where each vertex v corresponds to a street intersection and e edge corresponds to a street segment. For each G , we rasterise it into a 256 x 256 pixel image as shown in Figure 5. We also calculate basic network statistics [6] such as

¹©2017 Google Inc.



Figure 2: The Greater London case study area boundary.



Figure 3: Examples of street level images from Google StreetView. ©2017 Google Inc.

network centrality and network density that are later used to test the learnt features of the images.

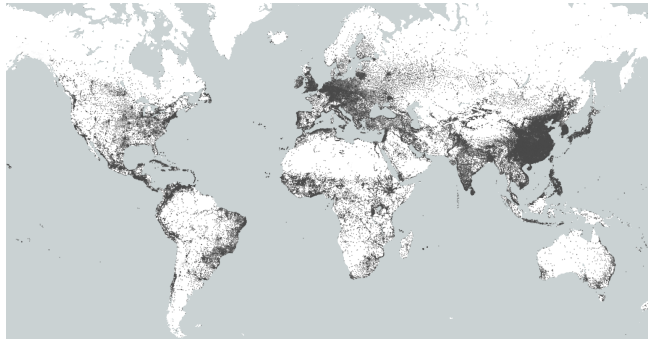


Figure 4: Centroid of 107,973 cities and towns used for training and testing.



Figure 5: Sample of rasterized street network data.

4 EXPERIMENTAL RESULTS

In order to discover new information and interpretations from these visual latent components, we will visualise these components and to use these factors for predictions on both the street level and street network dataset.

4.1 streetview images

4.1.1 Visualisation experiments. The *ConvPCA* first learns a mapping from a three channel street level images ($224 \times 224 \times 3$) down to a lower dimensional embedding (4,096 dimensions) using a convolutional autoencoder *CAE*. The lower dimension embedding is then further summarised into a set of uncorrelated components using PCA_{lin} . For the StreetView images, we adopted a *VGG* [30] like architecture where we keep the kernel size and filter numbers constant across both the encoder and decoder.

To show the results, we first plot the images with the highest and lowest principal component values for interpretation. In this case, component *pca 7* has blank facade in one of the extremes and natural scenarios in the other. *pca 10* and *pca 30* shows a tunnel space in one extreme and a mixture of urban scenarios in the other. While *pca 14* has buildings in one extreme and blank facades in the other. Lower rank components that capture lower variance seem to be showing less patterns and therefore not visualised.

To interrogate the results of the primary components, we focus on visualising *pca 1* and *pca 3* geographically in figure 6. The images plotted above the map show the two extremes of the visual latent components. We can descriptively interpret these two visual components v_1 and v_3 as proxy measures for different type of street urbanity. We also show through global spatial autocorrelation analysis these components exhibit strong spatial dependence. Please see the *appendix* for more details on the spatial analysis.

We then visualised one of the StreetView images and perturbed each of the two principal components while holding all the other component values constant before passing it to the decoder to generate a synthetic image. More formally, for each *pca* we create a mean vector \hat{v} , where we keep all values in \hat{v} constant and vary only the individual *pca* before passing it to the decoder to generate a synthetic image. Figure 9 shows when we perturbed *pca 1* of a typical StreetView image while holding all other principal components constant, building details tends to increase, and when we perturbed the same image in the other axis, building details tend to be reduced. In contrast, when we perturbed *pca 3* of the same StreetView image whilst holding all other principal components constant, trees started appearing and when we perturbed the same image to the other axis buildings becomes more prominent and the trees disappeared. The result also shows that the streets are widening in one of the axis while the car is disappearing in the other axis for *pca 1*. This result suggests, each component is related to a quality measure of street urbanity and is possibly capturing multiple correlated visual features of a StreetView image. As a result, in terms of controllability, the approach seems not able to disentangle highly correlated features. These descriptive results show geographical and generative visualisations are useful approaches to discover meanings from these visual latent components. However, more researches is needed to validate the meaning of these visual components quantitatively.

4.1.2 Prediction experiment. In order to demonstrate the usefulness of these visual latent components for different down-stream tasks, we constructed two separate models to map the latent visual components V to street enclosure (regression task) and street frontage quality (classification task). We compare PCA_{lin} to two

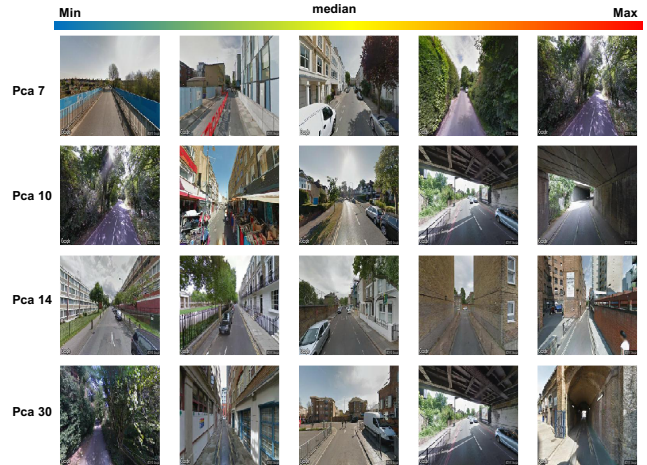


Figure 6: Example London Google StreetView images according to its component values. *pca 7* has blank facade in one of the extremes and natural scenarios in the other. *pca 10* and *pca 30* shows a tunnel space in one extreme and a mixture of urban scenarios in the other. *pca 14* shows buildings in one extreme and blank facades in the other. ©2017 Google Inc.

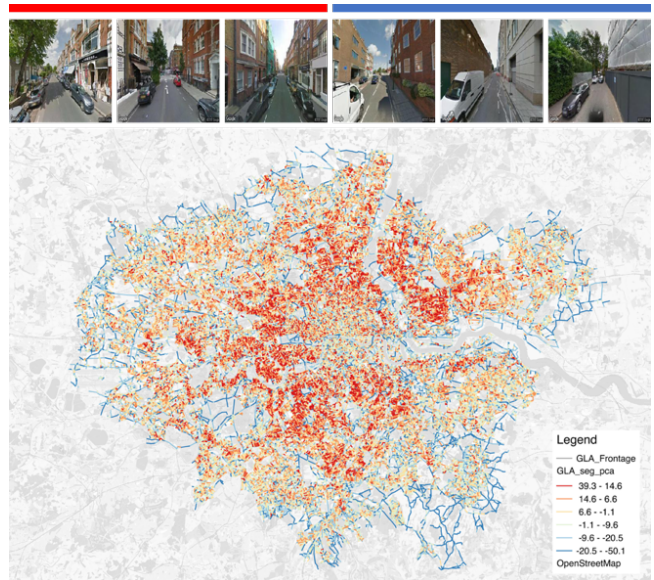


Figure 7: *pca 1* shows descriptively a quality measure for street urbanity. The map shows that the intra-urban area has higher component value. Visualising the extremes of the component shows greater building details in one end of the axis. ©2017 Google Inc.

other dimension reduction techniques in retrieving the latent visual component, a linear autoencoder AE_{lin} and a nonlinear autoencoder AE_{non} . The linear autoencoder uses a linear activation function with one bottleneck layer that outputs V . The nonlinear autoencoder on the other hand uses the *ReLU* activation function with three hidden layers where the first and the third hidden layer



Figure 8: *pca 3* shows descriptively another quality measure for street urbanity. The map shows central London has higher component values and outer London has less. Visualising the extremes of the latent component shows higher building density in one end of the axis. ©2017 Google Inc.



Figure 9: Visual latent component perturbations of a typical StreetView image. The first row shows the perturbation of *pca 1* where building details increase in one end of the axis and reduces on the other. The second row shows the perturbation of *pca 3* where greenery increase on one end of the axis and building density on the other end of the axis. ©2017 Google Inc.

are the encoding and decoding layer with 512 neurons and the second layer being the bottleneck layer that outputs V .

Street enclosure here is defined as the average height of the building of a street divided by the average width between the buildings of the same street as illustrated in fig 10. The street enclosure is calculated by segmenting the streets from Ordnance Survey data [32] every 40m. For each street segment S , we calculate the geographic median S_c and the azimuth S_α , and create a new line that is perpendicular to S at the point S_c . The perpendicular line S_\perp is used to create the street profile by intersecting it with the closest building on either side of the street and querying the associated height attribute, this is used to calculate the street enclosure as

building height to street width ratio $enc = \bar{h}/w$. Please see [25] for additional details.

Street frontage types here is defined with four frontage categories namely, active frontage on both sides of the street, active frontage on one side of the street, non-active frontage and non urban frontage. This dataset was manually compiled and studied from a previous study. Please see [19] for additional details.

We split the dataset randomly into a train (70%), validation (15%) and test set (15%). We then train a multi-layer perceptron $F(\cdot)$ to predict street enclosure and street frontage types from the visual latent components V as inputs, parameterized by a set of weights W_v .

The multi-layer perceptron (*mlp*) here is defined as a fully connected neural network with three hidden layers. The first fully connected layer has 64 hidden nodes, while the second layer has 32 hidden nodes and the third layer has 16 hidden nodes. A dropout layer (0.2) and a $l1$ regularisation was added in the final activation layer for better generalisation. To test the importance of the visual components with respect to the model accuracy, we constructed five different models based on the number of components [4,8,16,32,64] using the three dimension reduction techniques. This results in 15 models in total.

We train the street enclosure model to minimize the mean squared error *mse* on a training set, using the ADAM [15] optimizer with an initial learning rate set at 0.001. We then report the mean squared error (*MSE*) and the coefficient of determination R^2 between the model prediction and the observed street enclosure for the spatially random test-set. Similarly, we train the street frontage model to minimise the categorical cross-entropy losses on the training set, using ADAM [15] optimiser with an initial learning rate set at 0.001. We then report the cross entropy losses and the accuracy which is simply the sum of correctly predicted frontage class over all samples. All the experiments are conducted with the Keras library [7] using a Tensorflow [1] back-end.

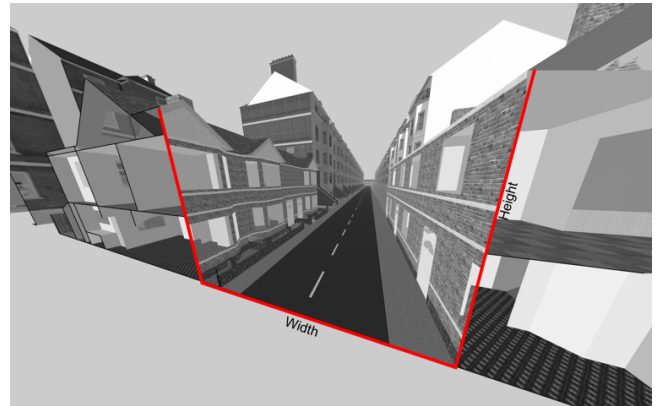


Figure 10: Street enclosure diagram. We define street enclosure as the ratio between $avg.height/avg.width$.

The results in table 1 shows the *losses* and *accuracy* of the three dimension reduction techniques when predicting street enclosure for a spatially random test-set. The model with 64 components achieve 60% accuracy, while the model with 4 components achieve 40-50% accuracy. The result shows, we can achieve similar levels of



Figure 11: Street frontage diagram. We classify street images into four street frontage categories namely *both-side-active*, *single-side-active*, *non-active*, *non-urban* [19]. ©2017 Google Inc.

accuracy with PCA_{lin} when compared to both AE_{lin} and AE_{non} . Similarly, the results in table 2 shows the *losses* and *accuracy* of the three dimension reduction techniques when predicting street frontage quality for a spatially random test-set. The results show a model with more components achieve a higher accuracy (70%) than one with less and that PCA_{lin} achieves comparable accuracy to AE_{lin} and AE_{non} . These results suggest, the convolutional layers are possibly capturing some of the non-linear effects between the different image features in the data. As a result, a linear dimension reduction technique such as PCA_{lin} , is able to learn a compact representation of the latent variable z which captures similar variance for two predictive tasks when compared to the autoencoders while retaining its interpretability.

Table 1: Street Enclosure Results

<i>accuracy</i>	PCA_{lin}	AE_{lin}	AE_{non}
4 components	41.50%	41.80%	50.96%
8 components	53.64%	54.02%	55.78%
16 components	56.37%	56.24%	56.77%
32 components	58.36%	58.00%	58.62%
64 components	59.17%	58.45%	59.98%

<i>losses</i>	PCA_{lin}	AE_{lin}	AE_{non}
4 components	0.603	0.600	0.499
8 components	0.480	0.474	0.450
16 components	0.450	0.451	0.440
32 components	0.429	0.433	0.421
64 components	0.421	0.428	0.407

4.2 street network

4.2.1 Visualisation experiments. For the street network case study, the trained convolutional autoencoder learned a mapping from the space of street network images ($256 \times 256 \times 1$ or 65,536 dimensions) to a lower dimensional latent space (640 dimensions) which are then further summarised into a set of linearly uncorrelated variables by applying (PCA_{lin}). By plotting out the street network images with the lowest to highest values of each component we can start to interpret the learnt latent space. In figure 12, we show the first five. These plots all relate to density of streets in different spatialised regions. The first pca encodes general density, while pca 2-5 encode spatialised densities (left-right, top-bottom, center-periphery, diagonals) respectively.

Table 2: Street Frontage Results

<i>accuracy</i>	PCA_{lin}	AE_{lin}	AE_{non}
4 components	61.46%	59.91%	60.96%
8 components	64.10%	62.26%	62.84%
16 components	68.24%	67.17%	67.44%
32 components	69.13%	68.51%	68.71%
64 components	71.41%	71.93%	70.50%

<i>losses</i>	PCA_{lin}	AE_{lin}	AE_{non}
4 components	0.884	0.898	0.907
8 components	0.839	0.846	0.858
16 components	0.758	0.759	0.772
32 components	0.734	0.735	0.751
64 components	0.707	0.669	0.719

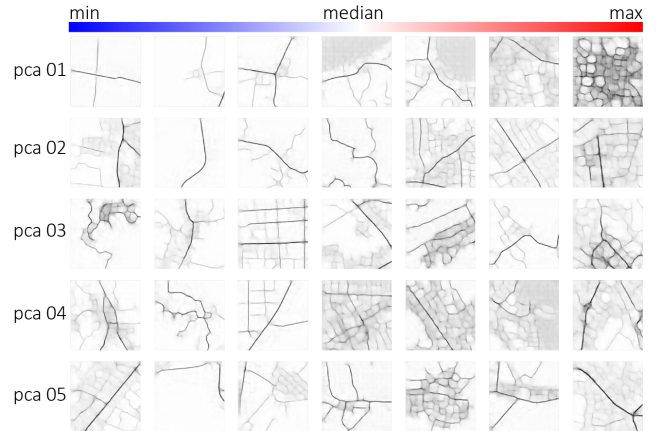


Figure 12: Example street network images for the first five principal components.

To make it easier to interpret each pca we create a mean vector \hat{v} , where we keep all values in \hat{v} constant and vary only the pca before passing it to the decoder to create a synthetic image. In figure 13, we show a subset of the different latent visual components encoded by the pca values. We show that the first 10 pca encode regions of spatialised density. We confirm the clustered spatial structure of these component through a spatial autocorrelation tests. The results of the test are shown in the *appendix* for the first 8 pca perturbations. pca 11-50 encode global structure of the network (coarse grain detail), while pca 50-640 encode local structure of the network (finer grain detail).

By mapping the values of the principal components we can further test spatial patterns that they might encode. With just the first principal component of the latent space we are able to differentiate street network densities across the city of London. Figure 14 shows central London has higher street density than outer London.

4.2.2 Prediction experiment. Lastly, we test the ability these encoding can capture network features by using them to predict two network statistics: intersection density and closeness centrality. To do so, we first select a number of cities from our dataset where we

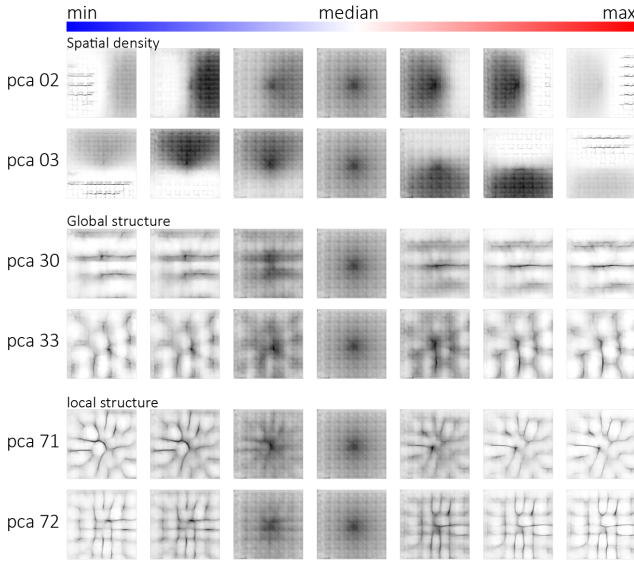


Figure 13: Latent visual component perturbation of an average street network image. By perturbing the visual component of an average image, we are able to show meaning of the perturbed component. The first sets of components seems to be related to spatialised density. While the second and third sets of components seems to be related to global and local structure of the street network.

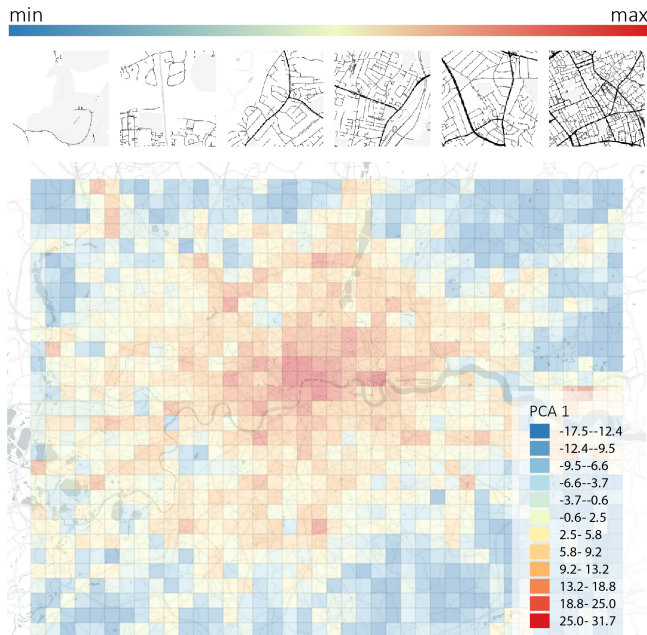


Figure 14: Values of the first principal component across London. The first component of the street network images can be interpreted as a measure of street network density.

retrieve its street network graphs $G = (V, E)$. For each graph, we calculate the closeness centrality of its nodes $u \in V$ through:

$$C(u) = \frac{n-1}{\sum_{v \in V} d(v, u)}$$

where $d(v, u)$ is the shortest weighted path between u and v and n is the total number of nodes in the graph. We then create a continuous $1.5 \times 1.5 \text{ km}$ rectangular grid over each city graph. For each grid cell, we define intersection density as the number of nodes inside the cell divided by the surface area, and the closeness centrality as its median values within each cell.

Table 3: Street intersection density results

	<i>accuracy</i>	PCA_{lin}	AE_{lin}	AE_{non}
4 components	76.59%	62.90%	73.56%	
8 components	77.28%	76.43%	72.01%	
16 components	75.54%	75.94%	74.00%	
32 components	71.15%	71.65%	76.00%	
64 components	69.45%	73.70%	71.83%	
	<i>losses</i>	PCA_{lin}	AE_{lin}	AE_{non}
4 components	0.23	0.37	0.25	
8 components	0.22	0.23	0.27	
16 components	0.26	0.23	0.26	
32 components	0.31	0.26	0.26	
64 components	0.35	0.24	0.25	

Table 4: Street closeness centrality results

	<i>accuracy</i>	PCA_{lin}	AE_{lin}	AE_{non}
4 components	54.63%	59.89%	60.23%	
8 components	55.41%	59.51%	58.22%	
16 components	54.17%	58.19%	59.02%	
32 components	51.85%	57.74%	58.92%	
64 components	34.33%	52.00%	53.45%	
	<i>losses</i>	PCA_{lin}	AE_{lin}	AE_{non}
4 components	0.46	0.43	0.43	
8 components	0.45	0.41	0.45	
16 components	0.48	0.46	0.43	
32 components	0.53	0.47	0.41	
64 components	0.61	0.48	0.47	

The data is then split into a train (70%), validation (15%), and test (15%) set. We train a $mlp F(\cdot)$ to predict both the intersection density and median closeness centrality for each grid cell for all our street network graphs from the visual latent components v . The mlp here is defined by a fully connected neural network with two hidden layers, and a dropout layer (0.2) before the final activation. We define five different models based on the number of components [4,8,16,32,64] using the three dimension reduction techniques described in the methods section.

The results in Table 3 show the mse and R^2 for street intersection density using different number of pca components. With just a few components we are able to achieve an accuracy of 77% with a

PCA_{lin} model for the spatially random test-set on spatial features of the graph (intersection density), achieving slightly better results than both AE_{lin} and AE_{non} .

In the case of the median closeness centrality, shown in Table 4, we achieve a R^2 of 55% with a PCA_{lin} , showing we can achieve similar levels of accuracy to the other dimensionality reduction techniques. The difference in results between the intersection density and the median closeness centrality predictions is most likely due to the fact the while the intersection density is a local attribute thus can be entirely captured through the local graph structure within the grid cell, closeness centrality is dependent on the global structure of the graph of the entire city. Despite this, the model is informative and is still able to capture a significant portion of the variance of the closeness value with few components of the local graph structure. Further research is needed to investigate how much of the global structure can be inferred by local attributes.

5 DISCUSSION AND CONCLUSION

We have presented a simple but novel unsupervised approach to extract and interrogate visual latent components from urban images. This exploratory research sits in contrast to previous works which focused on supervised learning [9, 19, 24, 29]. Through geographic mapping, generative visualisations, and prediction experiments we were able to retrieve initial meaning from these visual latent components. With the increasing availability of large scale unlabelled image data, research into learning a compact representation automatically from geographical data will become increasingly useful.

In the case of the street level images, by mapping the visual latent components and generating synthetic images by perturbing its components, we were able to discover descriptive meaning from the data of which two of the primary components could be measures in describing street urbanity. We also found the lower dimensional latent components are able to predict two different generic urban characteristics such as street enclosure and street frontage type. The predictive accuracy for street frontage type is not as high as those using a purely supervised approach [19] but the results suggest a useful and generalisable representation can be learnt for different tasks. Despite the positive results, further exploration is necessary. For example, research is needed to relate the principal components to humanly labelled data describing the perception of street quality [24]. The results can validate the meaning of these components quantitatively. To confirm the usefulness of the representation learnt, further research is also needed in comparing the visual latent components from unsupervised visual features with the visual latent components from supervised visual features (ie. Places365 database [34]). Future researches are also needed on a) creating more realistic reconstructions by using generative models such as VAE or GAN b) developing quantitative methods to systematically disentangle and control interpretable latent components and c) conducting future research and designing experiments on semi-supervised learning and multi-task learning tasks.

In the case of the street networks, although the model is able to predict road network density and median closeness centrality, it fails to capture more complex street network features, we believe this is because the self-organized pattern of street networks

is the result of both geometrical order/disorder as well as local rules of optimality. Through rasterising the street networks, the explicit topological data of the graph is lost, and the model is not able to recover this quality from the image alone. Future works can explore ways to incorporate topological properties of the networks into the model. Recent advances in graph neural networks provide promising directions that would allow both topological and geometric properties to be incorporated into the model, this would allow a richer representation of the street network as both local connectivity structure and their spatial embedding could be preserved. Despite its many limitations, there are benefits to such an approach where traditional network measures can sometimes be computationally expensive, for example *betweennesscentrality* has a time complexity of $O(nm + n^2 \log n)$ and many spectral properties require eigenvalue decomposition of the graph laplacian matrix to be computed, with a time complexity of $O(n^3)$. A model that could approximate these parameters in an efficient manner could prove useful for varied applications, such as characterizing street networks across the world.

An immediate implication of the study, is that by learning a useful and compact representation from urban images, we can use this information immediately for other down-stream geographical tasks such as in prediction and classification. Conversely, this can reduce compute time and data collection costs significantly. More importantly though, the exploratory knowledge discovery process of using a combination of visualisation and inference, can shed new information about these non-linear methods such as neural networks and higher dimensional complex datasets such as images. To conclude, this research contributes to recent efforts in linking the disciplines of geography and machine learning. On the one hand, we find meaning from the visual latent components of street level and street network images. On the other hand, we also demonstrate how geographical datasets and visualisation techniques can be useful to enrich our understanding of machine learning methods.

REFERENCES

- [1] Martin Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2015. TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems. <https://www.tensorflow.org/> Software available from tensorflow.org.
- [2] Michael Batty. 2005. Cities and complexity: understanding cities through cellular automata, agent-based models and fractals.
- [3] Yoshua Bengio, Pascal Lamblin, Dan Popovici, and Hugo Larochelle. 2007. Greedy Layer-Wise Training of Deep Networks. In *Advances in Neural Information Processing Systems 19*, B. Schölkopf, J. C. Platt, and T. Hoffman (Eds.). MIT Press, 153–160. <http://papers.nips.cc/paper/3048-greedy-layer-wise-training-of-deep-networks.pdf>
- [4] Geoff Boeing. 2017. OSMnx: New methods for acquiring, constructing, analyzing, and visualizing complex street networks. *Computers, Environment and Urban Systems* 65 (2017), 126–139.
- [5] Geoff Boeing. 2018. Measuring the Complexity of Urban Form and Design. October (2018), 1–22.
- [6] Geoff Boeing. 2018. A multi-scale analysis of 27,000 urban street networks: Every US city, town, urbanized area, and Zillow neighborhood. *Environment and Planning B: Urban Analytics and City Science* (2018), 2399808318784595.
- [7] François Chollet. 2015. keras. <https://github.com/fchollet/keras>.
- [8] Abhimanyu Dubey, Nikhil Naik, Devi Parikh, Ramesh Raskar, and Cesar Augusto Hidalgo. 2016. Deep Learning the City : Quantifying Urban Perception At A

- Global Scale. *European Conference on Computer Vision (ECCV)* (2016).
- [9] Timnit Gebru, Jonathan Krause, Yilu Wang, Duyun Chen, Jia Deng, Erez Lieberman Aiden, and Fei-Fei Li. 2017. Using deep learning and Google Street View to estimate the demographic makeup of neighbourhoods across the United States. *PNAS* (2017).
- [10] Google. 2018. <https://www.maps.google.com/>.
- [11] Xifeng Guo, Xinwang Liu, En Zhu, and Jianping Yin. 2017. Deep Clustering with Convolutional Autoencoders. *Neural Information Processing. ICONIP 2017. Lecture Notes in Computer Science* (2017).
- [12] Stefan Hartmann, Reinhard Klein, Raoul Wessel, and Reinhard Klein. 2017. StreetGAN : Towards Road Network Synthesis with Generative Adversarial Networks. (2017).
- [13] Bill Hillier. 2007. *Space is the machine: a configurational theory of architecture*. Space Syntax.
- [14] Kira Kempinska and Murcio Roberto. 2019. Modelling urban networks using Variational Autoencoders. *arXiv preprint* (2019). arXiv:arXiv:1905.06465v1
- [15] Diederik P. Kingma and Jimmy Ba. 2014. Adam: A Method for Stochastic Optimization. arXiv:cs.LG/1412.6980
- [16] Konstantin Klemmer, Adriano Soares Koshiyama, and Sebastian Flennerhag. 2019. Augmenting correlation structures in spatial data using deep generative models. *ArXiv abs/1905.09796* (2019).
- [17] Saïd Ladjal, Alasdair Newson, and Chi-Hieu Pham. 2019. A PCA-like Autoencoder. arXiv:cs.CV/1904.01277
- [18] Stephen Law, Brooks Paige, and Chris Russell. 2019. Take a Look Around: Using Street View and Satellite Images to Estimate House Prices. *ACM Transaction Intelligent Systems and Technology* 10, 5 (2019).
- [19] Stephen Law, Chanuki Illushka Seresinhe, Yao Shen, and Mario Gutierrez-Roig. 2018. Street-Frontage-Net: urban image classification using deep convolutional neural networks. *International Journal of Geographical Information Science* 0, 0 (2018), 1–27. <https://doi.org/10.1080/13658816.2018.1555832> arXiv:https://doi.org/10.1080/13658816.2018.1555832
- [20] Rémi Louf and Marc Barthélemy. 2014. A typology of street patterns. *Journal of The Royal Society Interface* 11, 101 (2014), 20140924.
- [21] Jonathan Masci, Ueli Meier, Dan Cireşan, and Jürgen Schmidhuber. 2011. Stacked Convolutional Auto-encoders for Hierarchical Feature Extraction. In *Proceedings of the 21th International Conference on Artificial Neural Networks - Volume Part I (ICANN'11)*. Springer-Verlag, Berlin, Heidelberg, 52–59. <http://dl.acm.org/citation.cfm?id=2029556.2029563>
- [22] Harvey J. Miller and Jiawei Han. 2001. *Geographic Data Mining and Knowledge Discovery*. Taylor & Francis, Inc., Bristol, PA, USA.
- [23] Vahid Moosavi. 2017. Urban morphology meets deep learning : Exploring urban forms in one million cities , town and villages across the planet. *arXiv preprint* (2017), 1–10. arXiv:arXiv:1709.02939v2
- [24] Nikhil Naik, Jade Philipoom, Ramesh Raskar, and Cesar Augusto Hidalgo. 2014. StreetScore - Predicting the Perceived Safety of One Million Streetscapes. In *CVPR Workshop on Web-scale Vision and Social Media*.
- [25] Mateo Neira and Laura Narvaez. 2019. The street as a three-dimensional urban form. In *ISUF 2019 XXVI international conference: Cities as Assemblages*.
- [26] OpenStreetMap contributors. 2019. Planet dump retrieved from <https://planet.osm.org> . <https://www.openstreetmap.org>.
- [27] Markus Reichstein, Gusta Camps-Valls, Bjorn Stevens, Martin Jung, Joachim Denzler, Nuno Carvalhais, and Prabhat. 2018. Deep learning and process understanding for data-driven Earth system science. *Nature* (2018). <http://www.nature.com/articles/s41586-019-0912-1>.
- [28] German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio M. Lopez. 2016. The SYNTHIA Dataset: A Large Collection of Synthetic Images for Semantic Segmentation of Urban Scenes. (2016). <refbase.cvc.uab.es/files/RSM2016.pdf>
- [29] C Seresinhe, T Preis, and S Moat. 2017. Using deep learning to quantify the beauty of outdoor places. *Royal Society Open Science* (2017).
- [30] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).
- [31] Emanuele Strano, Vincenzo Nicosia, Vito Latora, Sergio Porta, and Marc Barthélemy. 2012. Elementary processes governing the evolution of road networks. *Scientific reports* 2 (2012), 296.
- [32] Ordnance Survey. 2017. <https://www.ordnancesurvey.co.uk/opendatadownload/products.html>.
- [33] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Nikolai Yakovenko, Andrew Tao, Jan Kautz, and Bryan Catanzaro. 2018. Video-to-Video Synthesis. In *NeurIPS*.
- [34] Bolei Zhou, Agata Lapedriza, Jianxiong Xiao, Antonio Torralba, and Aude Oliva. 2014. Learning Deep Features for Scene Recognition using Places Database. In *Advances in Neural Information Processing Systems 27*, Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger (Eds.), Curran Associates, Inc., 487–495. <http://papers.nips.cc/paper/5349-learning-deep-features-for-scene-recognition-using-places-database.pdf>

6 APPENDIX

In the appendix, we describe for both the StreetView case study and the Street Network case study, the architecture of the Convolutional Autoencoder, the stacked autoencoders, the multi-layer-perceptron (*mlp*) and the spatial autocorrelation analysis.

6.1 StreetView architecture

6.1.1 Convolutional Autoencoder architecture. For the Convolutional Autoencoder of StreetView, the input is a fixed sized 224×224 three channel coloured image. We adopted a simplified convolution blocks from VGG [30] as the basis of the architecture where we keep the kernel size and filter numbers constant across both the encoder and decoder. Let Ck denote a Convolution-ReLU layer with k filters and Cdk denotes a Convolutional-ReLU-Upsample layer with k filters. All convolutions are 33 spatial filters applied with stride of 1.

encoder:C64-C64-C128-C128-C256-C256-C512-C512-CC512-C512-C512-CC512

decoder:CD512-CD512-CD512-C512-C512-C512-C256-C256-C256-C128-C128-C64-C64

6.1.2 Stacked Autoencoder architecture. For the stacked Autoencoder, where we summarise the latent variable z to its latent component v , we applied two forms of the autoencoder namely a linear autoencoder and a nonlinear autoencoder. The linear autoencoder uses linear activation functions with one bottleneck layer that outputs V . The nonlinear autoencoder on the otherhand uses the *ReLU* activation function with three hidden layers where the first and the third hidden layer are the encoding and decoding layer with 512 neurons and the second layer being the bottleneck layer that outputs V . Let Dk denote a Dense-ReLU layer with k filters and N as the number of components in the bottleneck layer.

linear: 4096-N-4096

non-linear: 4096-D512-N-D512-4096

6.1.3 Multi-layer Perceptron. The multi-layer perceptron *mlp* here is defined as a fully connected neural network with three hidden layers. The first fully connected layer has 64 hidden nodes, the second has 32 hidden nodes, while the third layer has 16 hidden nodes. A dropout layer (0.2) and l_1 regularisation was added in the final activation layer. We constructed five different models based on the number of components [4,8,16,32,64] and based on the three dimension reduction techniques resulting in a total of 15 models. Let Dk denote a Dense-ReLU layer with k number of neurons and N denote the shape of the visual latent component. [4,8,16,32,64].

Multi-layer perceptron: N-D64-D32-D16-1

6.2 Street network architecture

6.2.1 Convolutional Autoencoder architecture. For the Convolutional Autoencoder of street network data, the input is a fixed sized 256×256 single channel gray-scale. We use a stack of convolutional-ReLU layers and transposed convolutional layers, with a fixed small receptive field: 3×3 and a convolution stride fixed to 2 pixels. Let Ck denote a Convolution-ReLU layer with k filters and TCk denotes a Transposed-Convolution-ReLU layer with k filters.

encoder:C15-C15-C15-C10-C10

decoder:TC10-TC10-TC15-TC15-TC1

6.2.2 *Stacked Autoencoder architecture.* For the stacked Autoencoder, where we summarise the latent variable z to its latent component v , we applied two forms of the autoencoder namely a linear autoencoder and a nonlinear autoencoder. The linear autoencoder uses linear activation functions with one bottleneck layer that outputs V . The nonlinear autoencoder on the otherhand uses the *ReLU* activation function with three hidden layers where the first and the third hidden layer are the encoding and decoding layer with 128 neurons and the second layer being the bottleneck layer that outputs V . Let Dk denote a Dense-ReLU layer with k filters and N as the number of components in the bottleneck layer.

linear: 640-N-640

non-linear: 640-D128-N-D128-640

6.2.3 *Multi-layer Perceptron.* The multi-layer perceptron *mlp* here is defined as a fully connected neural network with two hidden layers. The first fully connected layer has 32 hidden nodes, while the second layer has 16 hidden nodes. A dropout layer (0.2) was added before the final activation layer and a $l1$ regularisation was added in the final activation layer. We constructed five different models based on the number of components [4,8,16,32,64] and based on the three dimension reduction techniques. Let Dk denote a Dense-ReLU layer with k number of neurons and N denote the shape of the visual latent component. [4,8,16,32,64].

Multi-layer perceptron: N-D32-D16-1

6.3 Global Spatial Autocorrelation structure

6.3.1 *Street network images.* In the case of the rasterized street network data, we test if the latent components capture strong local spatial inter-dependencies. This can be examined by measuring the autocorrelation of pixels with its local neighbours when perturbing the principal components of an average image. For our purposes we assume that the output of our *ConvPCA I'* follow some spatial process $y \sim f(c)$, where $y = \text{vec}(I')$ and c is a vector indexing the y_i pixel values in the output I' . The local spatial autocorrelation $L_i = L(y_i)$ is computed as:

$$L_i = (n - 1) \frac{y_i - \bar{y}}{\sum_{j=1, j \neq i} (y_j - \bar{y})^2} \sum_{j=1, j \neq i} w_{i,j} (y_j - \bar{y})$$

where \bar{y} represents the mean of y_i 's and $w_{i,j}$ are components of a weight matrix indicating membership of the local neighbourhood set between pixels i and j . Below, we show the results of the global spatial autocorrelation $L = \sum_i L_i$ of the max and minimum perturbations of the first 8 *pca* values and their corresponding I' s.

Table 5: Global spatial autocorrelation of outputs of min-max perturbations of first 8 PCA's of street networks

PCA	<i>Iofmin</i>	<i>Iofmax</i>
1 st component	0.87	0.88
2 nd component	0.96	0.87
3 rd component	0.93	0.89
4 th component	0.98	0.94
5 th component	0.99	0.97
6 th component	0.98	0.96
7 th component	0.99	0.98
8 th component	0.99	0.98

6.3.2 *StreetView images.* In the case of the street level images, we test if the latent components exhibit strong geographical dependencies. This can be examined by measuring the spatial autocorrelation between a street component value with its local neighbours, in this case defined by its 8th nearest local neighbours. The global Moran's IL for the two primary components are calculated. The result shows a strong spatial dependence of the visual latent component values at the street level.

Table 6: Global spatial autocorrelation of *pca* components

PCA	L
1 st component	0.68
3 rd component	0.75