

A multi-view contrastive learning framework for spatial embeddings in risk modelling

Freek Holvoet^{1,2}, Christopher Blier-Wong³, and Katrien Antonio^{1,2,4}

¹*Faculty of Economics and Business, KU Leuven, Belgium.*

⁴*Faculty of Economics and Business, University of Amsterdam, The Netherlands.*

²*LRisk, Leuven Research Center on Insurance and Financial Risk Analysis, KU Leuven, Belgium.*

³*Department of Statistical Sciences, University of Toronto, Canada*

November 25, 2025

Abstract

Incorporating spatial information, particularly those influenced by climate, weather, and demographic factors, is crucial for improving underwriting precision and enhancing risk management in insurance. However, spatial data are often unstructured, high-dimensional, and difficult to integrate into predictive models. Embedding methods are needed to convert spatial data into meaningful representations for modelling tasks. We propose a novel multi-view contrastive learning framework for generating spatial embeddings that combine information from multiple spatial data sources. To train the model, we construct a spatial dataset that merges satellite imagery and OpenStreetMap features across Europe. The framework aligns these spatial views with coordinate-based encodings, producing low-dimensional embeddings that capture both spatial structure and contextual similarity. Once trained, the model generates embeddings directly from latitude-longitude pairs, enabling any dataset with coordinates to be enriched with meaningful spatial features without requiring access to the original spatial inputs. In a case study on French real estate prices, we compare models trained on raw coordinates against those using our spatial embeddings as inputs. The embeddings consistently improve predictive accuracy across generalised linear, additive, and boosting models, while providing interpretable spatial effects and demonstrating transferability to unseen regions.

Keywords: spatial embedding, multi-view learning, GeoAI, spatial representation learning, risk analysis, satellite imagery

1 Introduction

Spatial data are central to fields such as environmental monitoring, urban planning, and real estate pricing. Spatial modelling is also increasingly important in insurance applications. For instance, spatial data are used to link weather patterns to hail-damage claim frequency (Gao and Shi, 2022) or to quantify the effect of heat waves on mortality rates (Robben et al., 2025). These examples illustrate how incorporating spatial information enables actuaries to better analyse

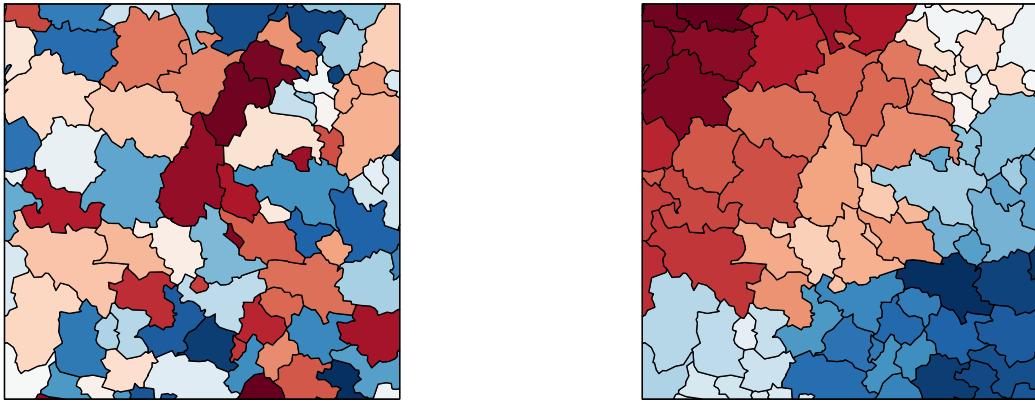
and manage risks that exhibit spatial dependence. Despite their importance, integrating spatial information into predictive models remains challenging. Location is often represented as a simple coordinate pair (latitude and longitude) or as a categorical variable, which fails to capture real-world geographical structure. Popular machine learning methods, such as random forests and gradient boosting machines, apply linear splits on these coordinate variables, leading to artificial rectangular patterns that do not align with real-world geography (Geerts et al., 2024).

Early work in spatial risk modelling was primarily inspired by geostatistics, where the focus is on understanding the dependence structures inherent in spatially distributed data. Two modelling strategies are typically explored: incorporating location via spatial smoothing of residuals, see Taylor (1989); Boskov and Verrall (1994); Taylor (2001) or by considering the spatial coordinates directly in a statistical model, for example, in a generalised additive model with a bivariate smooth effect of latitude and longitude, see Dimakos and Di Rattalma (2002); Gschlößl and Czado (2007); Shi and Shi (2017); Henckaerts et al. (2018, 2021). While these approaches model spatial dependence explicitly, modern machine learning increasingly represents spatial information through representation learning, where spatial information is captured implicitly via so-called embeddings.

Embeddings are vector representations of data objects, obtained through an encoder that maps discrete or structured inputs into a continuous, typically low-dimensional space. The dimensionality of the embedding determines how much information it can represent: low-dimensional embeddings capture the most essential structure of the input data, while higher-dimensional embeddings can preserve more detailed or complex patterns. A particular class of embeddings, called entity embeddings, refers to compact, continuous-valued representations of categorical variables, where each category is mapped to a low-dimensional vector (Guo and Berkhahn, 2016). The use of embeddings in actuarial science is not new: Richman (2021) discuss the notion of supervised entity embeddings to model categorical variables, and Blier-Wong et al. (2021) construct entity embeddings via unsupervised learning for use in actuarial applications. For example, a motor third-party liability dataset may contain categorical variables such as gender, vehicle type, and coverage level, each of which can be represented by entity embeddings. A non-exhaustive list of applications of entity embeddings in non-life insurance pricing includes Delong and Kozak (2023); Avanzi et al. (2024); Richman and Wüthrich (2024); Wilsens et al. (2024); Holvoet et al. (2025). Spatial information can be categorical in nature, for instance, when postal codes or administrative regions are treated as levels of a categorical variable. In such cases, entity embeddings can be used to capture spatial variation in risk levels, as shown in Shi and Shi (2023). In Blier-Wong et al. (2022), the authors claim that entity embeddings are inappropriate for integrating spatial information in insurance loss modelling. Granular territories such as postal codes are typically created to facilitate mail delivery, not to capture homogeneous risk profiles. Furthermore, if neighbors live on the border between two territories, their insurance premiums may differ since the embeddings of their territories will vary, even though their spatial risk should be similar. For this reason, the effect of spatial information in prediction models should be smooth on a map, and entity embeddings typically do not possess this property.

For these reasons, Blier-Wong et al. (2022) propose modelling spatial risk through spatial embeddings rather than entity embeddings. Spatial embeddings represent geographic information directly in a continuous latent space, where locations that are close in space are assigned similar vectors. Unlike entity embeddings, which treat each spatial unit as an independent category, spatial embeddings encode smooth transitions across space. Figure 1 illustrates the difference between entity embedding and spatial embedding, by embedding postal codes into one dimension, which is shown

as a colour on the map. In the left panel, postal codes treated as categorical entities produce random, discontinuous colour patterns, while the spatial embeddings in the right panel exhibit spatial consistency, with nearby postal codes having similar colours and distant postal codes having dissimilar colours.



- (a) Entity embeddings: each postal code is treated as a distinct category, resulting in random, discontinuous colour patterns that lack spatial structure.
- (b) Spatial embeddings: the embedding encodes geographic proximity, producing smooth colour transitions across neighbouring postal codes with nearby locations having similar embeddings, while distant locations differ.

Figure 1: Conceptual illustration of entity versus spatial embeddings for a selection of postal codes in France. Each postal code is embedded into a one-dimensional latent space, illustrated by colour.

Spatial embeddings can also incorporate external contextual information about each location, such as census statistics, land cover maps, or satellite imagery. These context-aware spatial embeddings enable models to represent not only the location itself, but also what is present there. In a case study on claim frequency modelling using Canadian data, [Blier-Wong et al. \(2022\)](#) construct spatial embeddings from geographic coordinates enriched with census data using a convolutional neural network. They show that these embeddings capture spatial risk patterns more effectively than traditional spatial smoothing methods, such as generalised additive models with bivariate smooth effects. Similarly, [Moussaid \(2023\)](#) construct spatial embeddings based on H3 hexagonal coordinates ([Uber Technologies Inc., nd](#)), augmenting them with spatial context from OpenStreetMap data for each hexagonal cell. Using a Belgian motor third-party liability claims dataset, they demonstrate that incorporating such contextual information improves predictive accuracy.

In the broader field of geospatial artificial intelligence (GeoAI), spatial representation learning has advanced rapidly through neural networks that process large-scale, heterogeneous spatial data ([Yin et al., 2019](#); [Blier-Wong et al., 2020](#); [Mai et al., 2020, 2022](#)). These methods typically produce task-agnostic spatial embeddings, which means that the embeddings are not trained for a single predictive task, but rather to capture general spatial structure and relationships across different regions or data sources. Once learnt, these embeddings can serve as general-purpose spatial representations that can be reused in different predictive settings. This property makes them particularly suited for spatial transfer learning, which refers to the ability of a model trained on one geographic area or dataset to generalise its learnt spatial patterns to other territories where no training data

are available. By leveraging shared geographic characteristics through spatial embeddings, transfer learning allows for more accurate predictions in previously unseen locations.

The recent development of foundation models extends context-aware spatial embedding models by aiming to create general-purpose geospatial representations that integrate multiple data modalities while retaining spatial awareness. Foundation models have already transformed language and vision through large pretrained architectures, and similar ambitions are now emerging for multimodal spatial models that unify imagery, vector data, and semantic context (Mai et al., 2024, 2025). A notable recent example is AlphaEarth (Brown et al., 2025), a large-scale multiview and spatiotemporal geospatial foundation model that integrates satellite imagery, weather reanalysis, and auxiliary contextual data to produce globally consistent embeddings. One specific spatial embedding architecture we will investigate in this paper is the SatCLIP model, proposed by Klemmer et al. (2023, 2025). SatCLIP constructs context-aware spatial embeddings by combining geographic coordinates with visual context extracted from satellite imagery. The model produces two types of embeddings: one derived from spatial coordinates using spherical harmonics to capture global spatial structure (Rußwurm et al., 2024), and another derived from satellite images using a vision transformer that extracts high-level visual patterns (Tarasiou et al., 2023). These two embedding spaces are aligned through contrastive learning, where the model learns to assign similar representations to spatial locations that share similar visual and geographic characteristics, and dissimilar representations to locations that differ. The contrastive training procedure allows the coordinate encoder to internalise the contextual information contained in the imagery. As a result, SatCLIP can generate embeddings for new locations using only their geographic coordinates, without requiring access to the original spatial data. Klemmer et al. (2025) uses the trained SatCLIP model to produce embeddings that are high-dimensional and designed to capture large-scale spatial variation across the earth’s surface. They demonstrate the use of these embeddings in a case study on global temperature pattern prediction.

In this paper, we propose a multi-view contrastive learning model that generates spatial embeddings by leveraging multiple input data sources as spatial context. Our contribution lies in combining information from different spatial views with geographical coordinates to obtain low-dimensional embeddings that capture local spatial patterns. The model takes as input different sources of spatial data, such as satellite imagery and tabular neighbourhood information, together with the geographical coordinates of where the spatial information was gathered. Through contrastive learning, the model learns to associate patterns in the geographic coordinates with the information contained in the spatial views. Rather than targeting global coverage, the proposed framework is intended for regional settings. It can be retrained efficiently on consumer-grade hardware, flexibly incorporates different spatial data sources, and yields fine-grained spatial embeddings that capture local geographic structure in greater detail when compared to global-scale spatial embedding models. Our work differs from SatCLIP (Klemmer et al., 2023, 2025) in two key aspects. First, we integrate multiple spatial data sources rather than relying on a single modality, such as satellite imagery. Second, our model is designed to learn compact embeddings optimised for fine-grained regional applications, whereas SatCLIP produces high-dimensional global representations tailored to large-scale earth system modelling. Our model architecture has two main components: one that encodes spatial views and another that encodes coordinates, as illustrated in Figure 2. The bottom of Figure 2 shows how each spatial input data source is processed by its dedicated encoder, represented by the blue blocks. The results of each spatial view encoder are then merged with a neural network part called the multi-view fusion encoder (orange block). The top half of Figure 2

shows the encoding of coordinates, represented by the green block in Figure 2. Hereto, we use a set of spherical harmonic basis functions (Rußwurm et al., 2024). These functions capture patterns on earth at both large and small scales. Each set of coordinates is encoded as a weighted sum of spherical harmonic basis functions. These weights are then passed through a neural network that learns to link the basis functions to meaningful geographic context. We use contrastive learning to connect the output of the multi-view fusion encoder and the output of the location encoder. This forces the neural network part of the location encoder to link the information from the multiple spatial views to the patterns from the spherical harmonic basis functions.

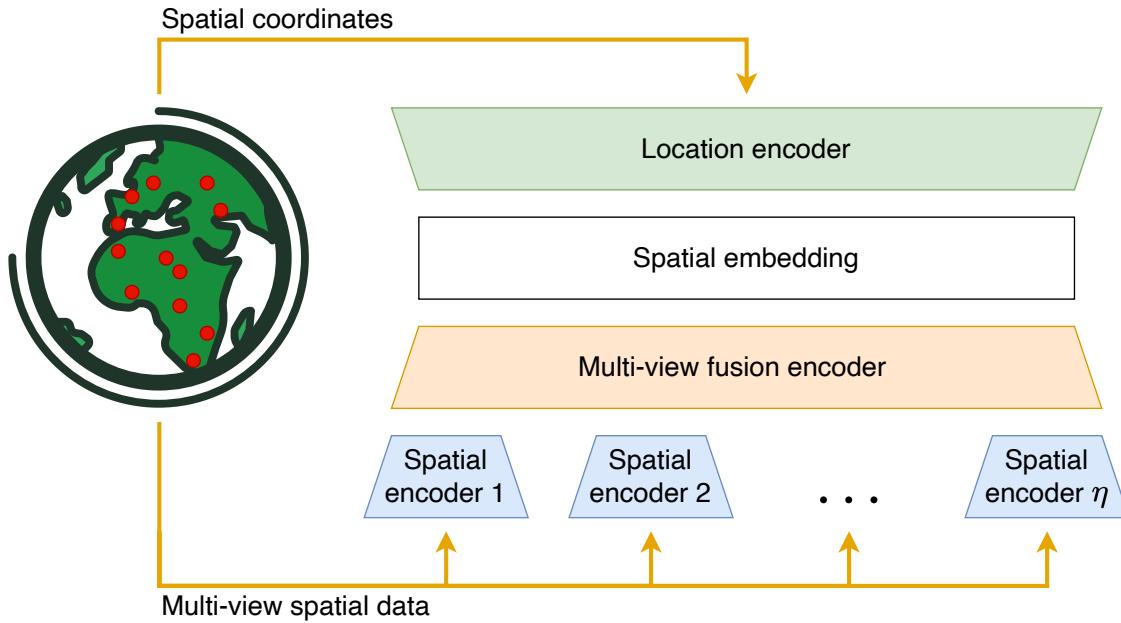


Figure 2: Architecture of our multi-view spatial embedding model. The green block represents the embedding of latitude and longitude coordinates. Each blue block embeds a different spatial view, and the orange block fuses these different spatial embeddings. The model is trained with a contrastive learning objective to learn spatial embeddings by contrasting the coordinate embedding with the information learnt from the different spatial views. Image inspired by Klemmer et al. (2025).

This paper is structured as follows. Section 2 details the construction of a multi-view spatial dataset for European countries. We gather both satellite images via Google Maps (Google, nd) and neighbourhood context via OpenStreetMap (OpenStreetMap, nd) for almost 100 000 locations across Europe. Section 3 shows the construction of the multi-view spatial embedding model and details the contrastive loss function, the training, and the hyperparameter choices. In Section 4, we illustrate the practical value of our spatial embeddings in the context of real estate pricing. Traditional models in this field often rely on basic location indicators, such as postal codes or proximity to landmarks, which only provide limited spatial context. In contrast, our proposed spatial embeddings incorporate a richer set of geographic features, enabling the model to better account for the influence of surrounding areas and environmental characteristics on property values. To evaluate our approach, we compare the predictive performance of a generalised linear model, a generalised additive model with a spatial smoother, and a gradient-boosting machine in predicting real estate prices using our spatial embeddings instead of raw coordinates. We examine permutation-based variable importance to evaluate the impact of spatial embeddings on the model’s output. Using

partial dependence effects, we demonstrate how geographic patterns, encoded via our embedding vectors, influence the predictions from the pricing model. We also experiment with transfer learning to areas held out of training to evaluate spatial generalization. Section 5 concludes the paper.

2 Dataset of locations and spatial views

We construct a multi-view dataset of 95 857 unique locations distributed across the European continent. The sampling of locations is guided by population distribution, using the European Local Administrative Units (LAU) framework (Eurostat, 2021). We show the distribution of the datasets’ locations in Europe in Figure 3.

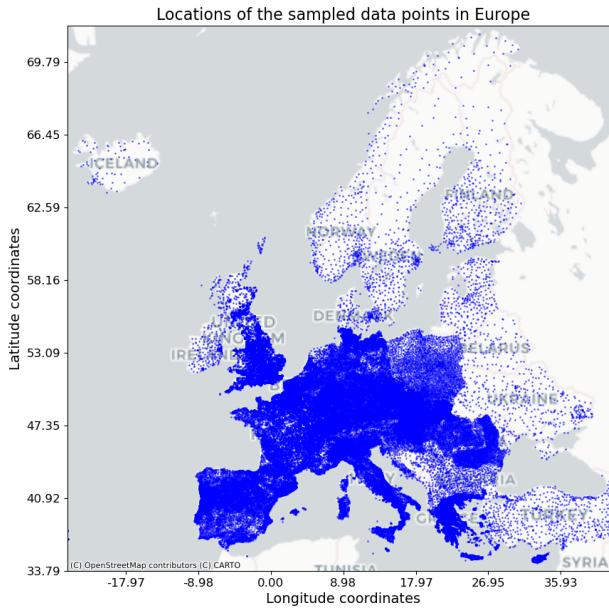


Figure 3: Our dataset consists of 95 857 sampled locations in European countries. The sampling of the locations is done based on the European Local Administrative Units framework.

Each location is defined by a latitude–longitude coordinate pair and enriched with two data views: a Google satellite image (GS) on the one hand and structured semantic context from OpenStreetMap (OSM) (OpenStreetMap, nd) on the other hand. Figure 4 shows an example of the two spatial views for a location in Leuven, Belgium. We gather the satellite image using the Google Maps Static API (Google, nd), as shown in Figure 4b. The satellite image is centred on the coordinates of the location and has a resolution of 256×256 pixels in RGB format, resulting in a dimension of $256 \times 256 \times 3$. We gather OSM tags representing semantic information typically not visible in satellite imagery. Hereto, we define six categories of tags: *shops*, *tourism*, *food & drink*, *health*, *education* and *public safety*. Appendix A contains all OpenStreetMap tags corresponding to each category. We aggregate occurrences of these tags in three rings of H3 hexagonal cells (Uber Technologies Inc., nd) centred at each location. Using concentric rings around the centre, we capture both close-by and more distant amenities in a consistent spatial structure. This is illustrated in Figure 4c, which shows three rings of hexagons around the centre location with the occurrences of each category of

tag in coloured points. Three rings of hexagons around a central location lead to 37 hexagons in total, where the frequency of each of the six tag categories is counted for each hexagon.

Formally, let $P = 95\,857$ denote the total number of locations in our dataset. We define the dataset $\mathcal{D} = \{(\lambda_p, \phi_p), V_p^{\text{GS}}, V_p^{\text{OSM}}\}_{p=1}^P$, where the pair (λ_p, ϕ_p) represents the geographical coordinates latitude and longitude of location p . Each location is associated with a satellite image represented by $V_p^{\text{GS}} \in \mathbb{R}^{256 \times 256 \times 3}$, and $V_p^{\text{OSM}} \in \mathbb{N}^{37 \times 6}$ represents the aggregated of OpenStreetMap tags counts. Note that for a location p , the spatial views V_p^{GS} and V_p^{OSM} always relate to the coordinates (λ_p, ϕ_p) .

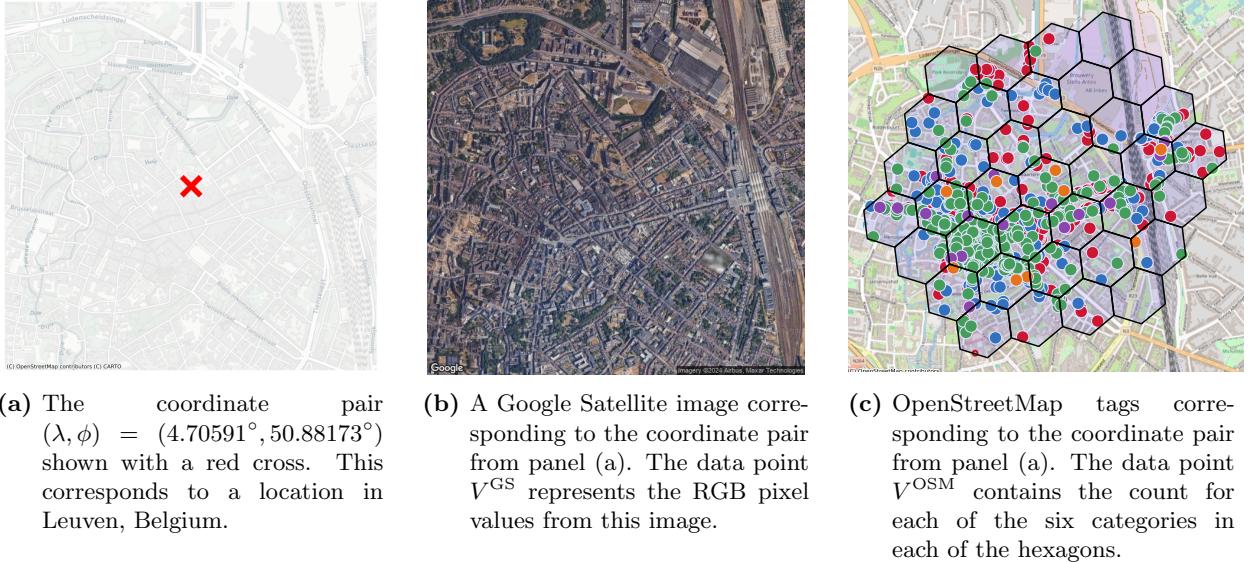


Figure 4: Illustration of one data point in the dataset \mathcal{D} , with coordinates $(\lambda, \phi) = (4.70591^\circ, 50.88173^\circ)$ in Leuven, Belgium. Panel (a) shows the location of the coordinate pair. Panel (b) displays the associated Google Satellite image, where the matrix $V^{\text{GS}} \in \mathbb{R}^{256 \times 256 \times 3}$ stores the RGB pixel values. Panel (c) displays the corresponding OpenStreetMap features, where the matrix $V^{\text{OSM}} \in \mathbb{N}^{37 \times 6}$ contains the counts of six categories of tags for each of the 37 surrounding hexagons.

3 Spatial embedding model architecture

Using the two spatial views introduced in Section 2, we construct our spatial embedding model. Firstly, we introduce the two spatial view encoders and the structure of the multi-view fusion encoder as sketched in Figure 2. Secondly, we detail the architecture of the location encoder. Thirdly, we provide an overview of the model’s complete architecture and examine the hyperparameter choices for training the spatial embedding model.

3.1 Spatial view encoders and multi-view fusion

Google satellite encoder (GS) The image encoder, denoted as $f^{\text{GS}}(V^{\text{GS}})$, maps a satellite image $V^{\text{GS}} \in \mathbb{R}^{256 \times 256 \times 3}$ to a d_{GS} -dimensional vector. As encoder f^{GS} , we use a convolutional neural network, since such architectures are specifically designed to extract and combine patterns

from image data. In particular, we adopt the ResNet18 model from Wang et al. (2023), a convolutional neural network trained on satellite imagery using momentum contrast (He et al., 2020), a self-supervised learning method that optimises the network to produce vector representations reflecting the structure of images. The ResNet18 architecture consists of 17 convolutional layers with residual links to facilitate training (He et al., 2016). Each convolutional layer consists of several convolutional filters, meant to extract patterns from the input data. After applying the convolutional layers, the results are pooled and flattened into a vector, which is then passed through a fully connected network to produce the final embedding. Figure 5 illustrates this architecture. The input satellite image V^{GS} is represented by a green input block. It is processed through a sequence of convolutional layers (blue blocks), followed by pooling and flattening operations (yellow blocks). The resulting representation is then passed through a fully connected network (rightmost blue block), which outputs the final embedding vector of dimension d_{GS} shown in red.

We initialise the model with pretrained weights from Wang et al. (2023) and adjust the output dimension to d_{GS} , a hyperparameter. The trainable parameters of f^{GS} lie in the convolutional filters and in the fully connected output network, represented as blue blocks in Figure 5.

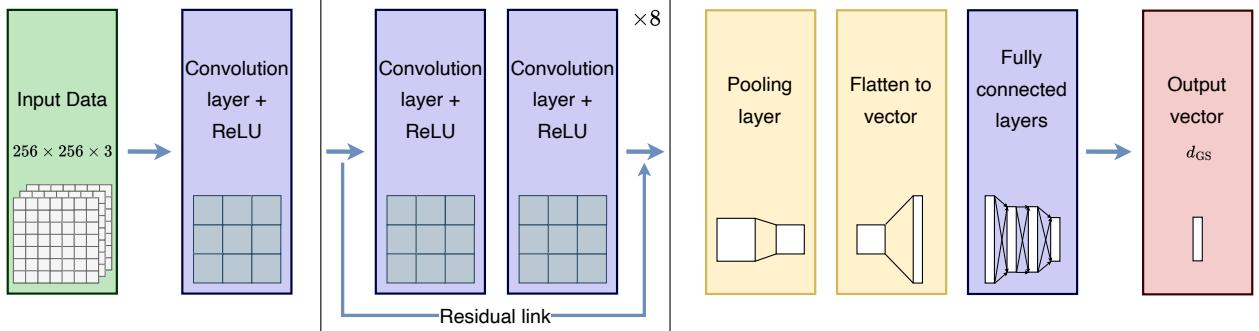
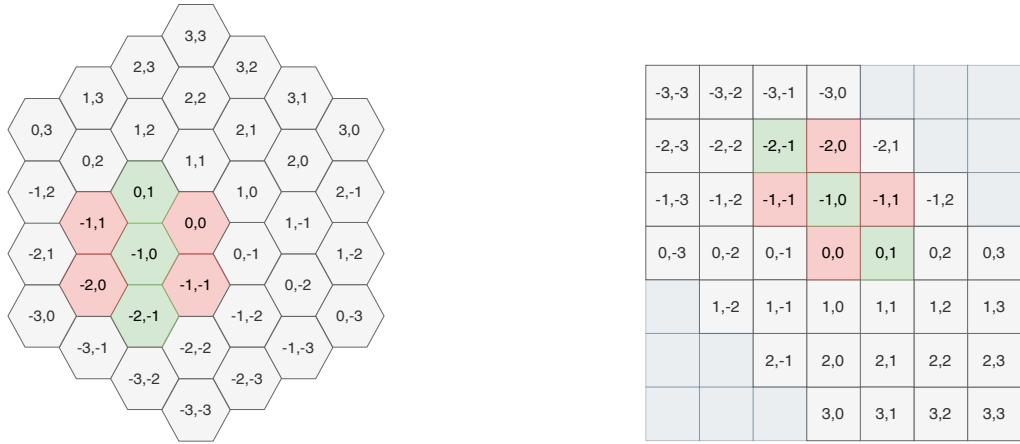


Figure 5: Illustration of the encoder used for the Google satellite imagery spatial view. We use a ResNet18 architecture (He et al., 2016), pretrained following Wang et al. (2023). The input image (green block) of dimension $256 \times 256 \times 3$ is processed by convolutional layers (blue blocks) that extract spatial patterns. These patterns are then pooled and flattened (yellow blocks) before passing through a fully connected network (rightmost blue block). The final embedding is a vector of dimension d_{GS} , as shown in the red block. Blue blocks contain the trainable parameters of the model.

OpenStreetMap tags encoder (OSM) We denote the encoder for the OpenStreetMap tags view as $f^{OSM}(V^{OSM})$ with the encoder mapping a hexagonal data input $V^{OSM} \in \mathbb{N}^{37 \times 6}$ to a d_{OSM} -dimensional vector. Similar to the encoder f^{GS} , we use a convolutional neural network to extract patterns from the input data and represent those patterns in a vector. However, our input is defined on a hexagonal grid, whereas standard convolutional neural networks are designed for square-grid data such as images. Therefore, we adopt the method proposed by Moussaid (2023) for convolutional neural networks on hexagonal input data while preserving neighbourhood relations between cells. The mapping follows the h2 coordinate system (Snyder et al., 1999), which minimises distortion by maintaining adjacency and approximate distance between neighbouring hexagons. Although alternative mappings are possible, this approach provides a consistent alignment between hexagonal neighbourhoods and square convolutional filters, allowing the use of standard convolutional layers without altering the underlying spatial structure of the data.

We illustrate this transformation with a grid of 37 hexagons in Figure 6a, which is mapped to a 7×7 matrix in Figure 6b following the h2 coordinate system from Snyder et al. (1999). The upper and lower corners of the matrix are masked, meaning that these entries are fixed to zero because no corresponding hexagons exist in those positions. For clarity, coordinates are added to each hexagonal cell and their correspondence to the square matrix is shown. A convolutional filter covering one ring of hexagons is illustrated with red and green cells in Figure 6a. In the transformed representation, this filter corresponds to a 3×3 convolutional filter, where the upper-right and lower-left corners are masked, shown with the red and green cells in Figure 6b. Generally, for a convolutional filter of k -rings around a central hexagon, the corresponding standard convolutional filter has a size of $2k+1$, with zeros masking the upper-right and lower-left corners of size k . Hence, the encoder f^{OSM} is built using a standard convolutional architecture on transformed input data and by fixing the upper right and lower left corners of the convolutional filters to zero.



(a) A hexagonal grid of 37 cells (three rings). Coordinates are added for clarity. In colour, a convolutional filter of size $k = 1$ (one ring) is shown.

(b) The corresponding 7×7 square grid obtained from the transformation of the 37 hexagonal cells following Moussaid (2023). Coordinates indicate how hexagonal cells map to the square representation. In colour, the $k = 1$ filter is shown after transformation, corresponding to a 3×3 filter with the upper-right and lower-left corners masked.

Figure 6: Illustration of the transformation from hexagonal to square representation. A convolutional filter of size $k = 1$ on the hexagonal grid (left) corresponds to a 3×3 filter on the transformed square grid (right), with the upper-right and lower-left corners masked.

The encoder f^{OSM} is illustrated in Figure 7. The input is a hexagonal grid consisting of 37 cells, where each cell counts the number of tags for each of the six OSM categories under consideration, resulting in an input dimension of 37×6 . We apply the transformation on our hexagonal input, which results in a square matrix of dimension $7 \times 7 \times 6$. On the transformed input, we apply two convolutional layers. The number of convolutional filters in each layer, together with the size of the filters, is treated as a tuning parameter. Each convolutional layer is followed by normalization and ReLU activation to facilitate training. The output of the second convolutional layer is then flattened into a vector and passed through a feed-forward neural network to produce the final representation of dimension d_{OSM} . This architecture enables the encoder to learn from the spatial structure of the OSM data, resulting in an output vector that accurately represents the spatial patterns in the input data.

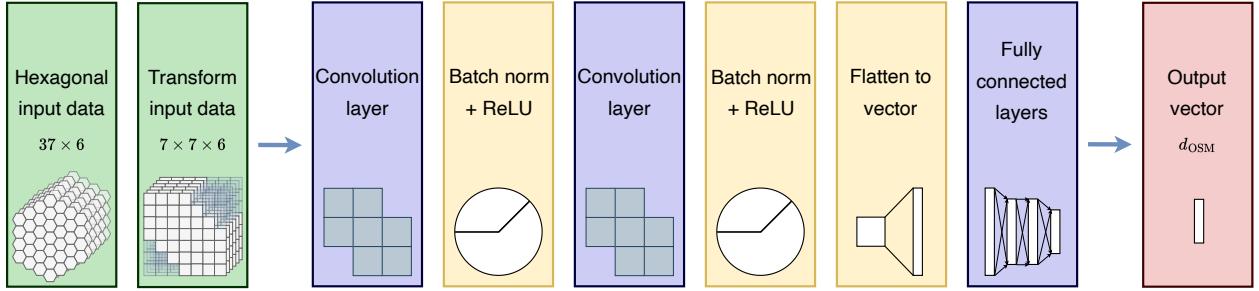


Figure 7: Illustration of the used encoder for the OpenStreetMap spatial view. The hexagonal input data is first transformed into a rectangular format. Next, there are two convolutional layers, each using masked convolutions. Finally, a feed-forward neural network outputs the encoded input. Each blue block in the flowchart contains trainable elements to be optimised in the training of the spatial encoder model.

The encoder f^{OSM} has trainable parameters in the filters of each of the two convolutional layers and in the fully connected network that produces the output. Figure 7 shows these components as blue blocks. The output dimension d_{OSM} is a hyperparameter. The number and size of convolutional filters in each layer are tuning parameters.

Multi-view fusion We combine the embedding vectors produced by the spatial view encoders f^{GS} and f^{OSM} into a single joint spatial representation by use of a multi-view fusion encoder f^{MV} . The two embeddings are concatenated and passed through a feed-forward neural network with fully connected layers, which outputs a vector of dimension d . These layers contain trainable weights that are optimised jointly with the rest of the spatial embedding model. The depth of the network and the width of each layer are treated as hyperparameters in the model design.

3.2 Location encoder

Each data point's coordinates (λ, ϕ) are embedded onto a d -dimensional vector by a location encoder, as illustrated in Figure 8. We denote the location encoder as

$$f^{\text{LE}} : \mathcal{D} \rightarrow \mathbb{R}^d : (\lambda, \phi) \mapsto f^{\text{LE}}(\lambda, \phi). \quad (1)$$

We use the proposed location encoder by Rußwurm et al. (2024), where $f^{\text{LE}}(\lambda, \phi) = \text{NN}(\text{PE}(\lambda, \phi))$, consisting of a positional encoder $\text{PE}(\lambda, \phi)$ and a neural network part $\text{NN}(\cdot)$. Rußwurm et al. (2024) apply this setup to learn spatial patterns from weather variables, i.e., temperature and air pressure data, for use in a regression model, and from satellite imagery for land-use classification and land-ocean segmentation.

The positional encoder is defined as

$$\text{PE}(\lambda, \phi) = \sum_{\ell=0}^L \sum_{m=-\ell}^{\ell} w_{\ell}^m Y_{\ell}^m(\lambda, \phi), \quad (2)$$

where w_ℓ^m are weights and Y_ℓ^m are the orthogonal spherical harmonic basis functions associated with Legendre polynomials P_ℓ^m of degrees ℓ and orders m , see Rußwurm et al. (2024). These spherical harmonic functions encode spatial locations on a spherical surface by decomposing them into a sum of orthogonal components. The orthogonality ensures that different spatial components are captured separately, allowing the model to represent complex spatial patterns at global and localised scales. The neural network part, known as the Siren network (Sitzmann et al., 2020; Klemmer et al., 2025; Rußwurm et al., 2024), takes as input the weighted spherical harmonic functions $w_\ell^m Y_\ell^m(\lambda, \phi)$. These are passed through feed-forward layers with sinusoidal activation functions, which allow the network to represent spatial patterns smoothly on a sphere.

Figure 8 sketches the architecture of the location encoder. The input is a coordinate pair (λ, ϕ) , represented by the green block. We evaluate the spherical harmonic basis functions (yellow blocks), for $\ell = 0, \dots, L$ and for each ℓ over $m = -\ell, \dots, \ell$, using this coordinate pair as input. The weighted spherical harmonic basis functions (first blue block) serve as inputs to the feed-forward neural network with sinusoidal activation functions, known as the Siren network (second blue block). The output of the Siren network is the location embedding of the coordinate pair.

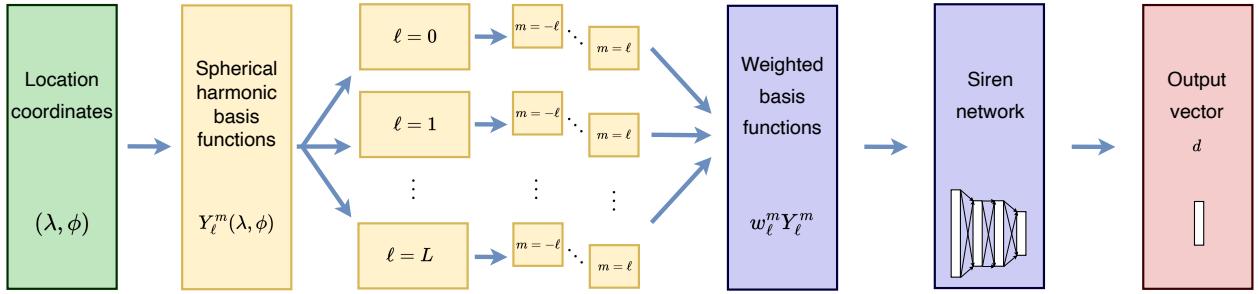


Figure 8: Illustration of the used location encoder. A weighted sum of spherical harmonic basis functions is used, with $\ell = 0, \dots, L$ and for each ℓ over $m = -\ell, \dots, \ell$. The weighted spherical harmonic basis functions $w_\ell^m Y_\ell^m$ are inputs to a feed-forward neural network using sinusoidal activation functions, called the Siren net (Rußwurm et al., 2024). Each blue block in the flowchart contains trainable elements that can be optimised during the training of the spatial encoder model.

The weights of the spherical harmonic basis functions and the weights in the Siren network are trainable parameters of the location encoder. Increasing the number of basis functions L allows the positional encoder to capture more fine-grained spatial detail, but also increases computational cost. To effectively utilise this additional detail, a larger value for L typically implies a larger value for d , ensuring the encoder has sufficient capacity to represent the information (Klemmer et al., 2025). The values of L , d , and the architecture of the Siren network are treated as hyperparameters.

3.3 Putting it all together: model architecture, training, and applications

Figure 9 illustrates the complete architecture for our spatial embedding model using two spatial views as inputs. The left-hand side represents the encoding of the location coordinates (λ, ϕ) by the location encoder f^{LE} as discussed in Section 3.2, resulting in an embedded location vector of dimension d . On the right, the spatial view encoders, as discussed in Section 3.1, process the two data sources: the OpenStreetMap tags, encoded by $f^{\text{OSM}}(V^{\text{OSM}})$, into a vector of dimension d_{OSM} , and the Google satellite imagery, encoded by $f^{\text{GS}}(V^{\text{GS}})$, into a vector of dimension d_{GS} .

sion d_{GS} . These two spatial view embeddings are then combined by the multi-view fusion encoder f^{MV} ($f^{GS}(V_p^{GS})$, $f^{OSM}(V_p^{OSM})$) into a single joint embedding of dimension d , with typically $d \leq d_{GS} + d_{OSM}$.

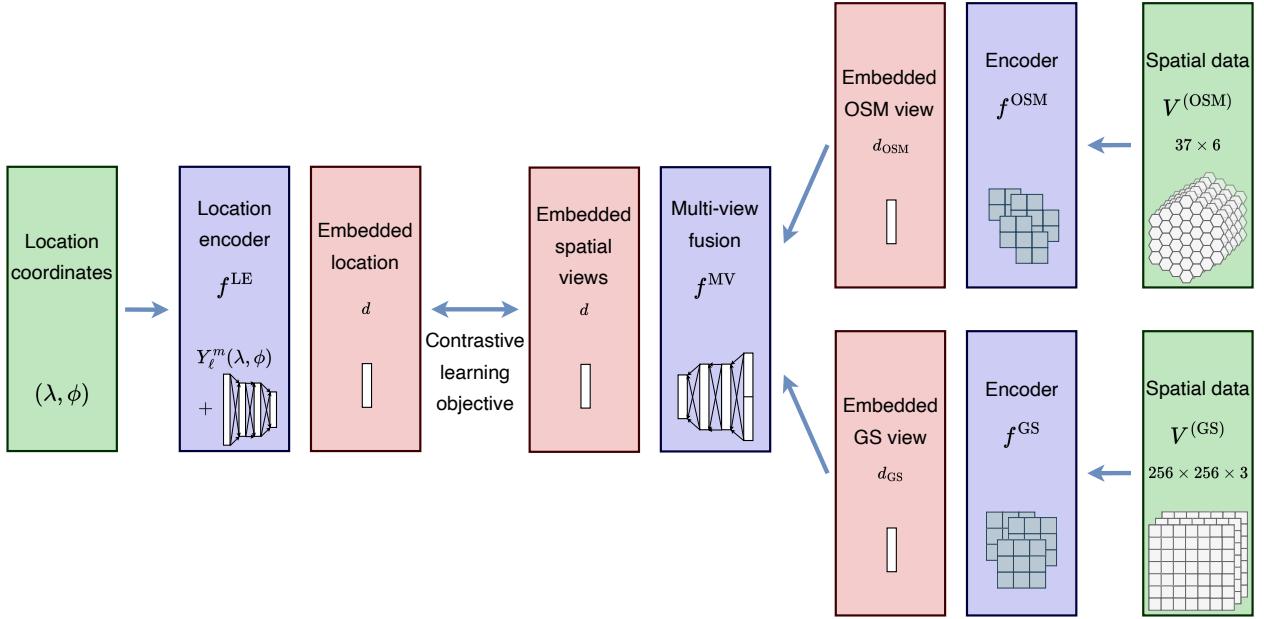


Figure 9: Architecture of the multi-view spatial embedding model with two spatial views as inputs. The left side shows the location encoder f^{LE} (Figure 8), which embeds the coordinates (λ, ϕ) . On the right, the Google satellite imagery and OpenStreetMap tags are processed by their respective encoders f^{GS} and f^{OSM} (Figures 5 and 7). Their outputs are combined in the multi-view fusion encoder f^{MV} . The resulting spatial view embedding is aligned with the location embedding using contrastive learning. Blue blocks indicate trainable components that are optimised during model training.

From the dataset \mathcal{D} , introduced in Section 2, we construct a training dataset $\mathcal{D}^T = \{(\lambda_p, \phi_p), V_p^{GS}, V_p^{OSM}\}_{p=1}^{P^T}$ of size P^T , containing 90% of the dataset \mathcal{D} sampled randomly. The remaining 10% of data forms the validation set \mathcal{D}^V and is used for early stopping in the training of the model. The training objective is to align the embedded locations with the embedded spatial views. For each $p = 1, \dots, P^T$, we denote the embedding vectors as z for notational clarity:

$$z_p^{LE} = f^{LE}(\lambda_p, \phi_p), \quad z_p^{MV} = f^{MV}(f^{GS}(V_p^{GS}), f^{OSM}(V_p^{OSM})).$$

The contrastive loss function \mathcal{L} is constructed to maximise the similarity between the so called positive pairs z_p^{LE} and z_p^{MV} , and to minimise the similarity across all mismatched pairs z_p^{LE} and $z_{p'}^{MV}$ for $p' \neq p$, see Radford et al. (2021). The contrastive loss over all elements in the training dataset is defined as

$$\mathcal{L} = \frac{1}{2P^T} \sum_{p=1}^{P^T} \left[-\log \frac{\exp(\text{sim}(z_p^{LE}, z_p^{MV})/\tau)}{\sum_{p'=1}^{P^T} \exp(\text{sim}(z_p^{LE}, z_{p'}^{MV})/\tau)} - \log \frac{\exp(\text{sim}(z_p^{MV}, z_p^{LE})/\tau)}{\sum_{p'=1}^{P^T} \exp(\text{sim}(z_p^{MV}, z_{p'}^{LE})/\tau)} \right],$$

where $\text{sim}(a, b) = \frac{a \cdot b}{\|a\| \|b\|}$ denotes the normalised dot-product and τ is a so-called temperature hyperparameter. The two similarity terms ensure that both the alignment from location to spatial

embeddings, $\text{sim}(z_p^{\text{LE}}, z_p^{\text{MV}})$, and from spatial to location embeddings, $\text{sim}(z_p^{\text{MV}}, z_p^{\text{LE}})$, is maximised. Training the complete architecture of our spatially embedded model consists of finding the optimal values for the trainable parameters so that the contrastive loss \mathcal{L} on the training dataset \mathcal{D}^T is minimised. The trainable parameters in the model architecture are illustrated with blue blocks in Figure 9. Minimizing the contrastive loss encourages the location embeddings z_p^{LE} and the multi-view embeddings z_p^{MV} of the same point p to be similar, while embeddings of different locations remain distinct. In this way, the positional patterns represented by the spherical harmonics are aligned with the information extracted from the spatial views. The resulting embedded locations thus encode both the geographic position and the spatial content, allowing the model to represent locations in a manner that reflects where they are as well as what is observed there.

In a case study involving geographical locations, the trained spatial embedding model can be applied to represent each location by its embedding. For a dataset consisting of covariates x_1, \dots, x_p , where x_1 represents geographical coordinates, the covariate x_1 can be replaced with its embedding $f^{\text{LE}}(x_1)$, as illustrated in Figure 10. These embeddings act as compact descriptors that integrate both geographic position and the spatial information captured from multiple data views, and can be used as input features in downstream analyses such as regression models.

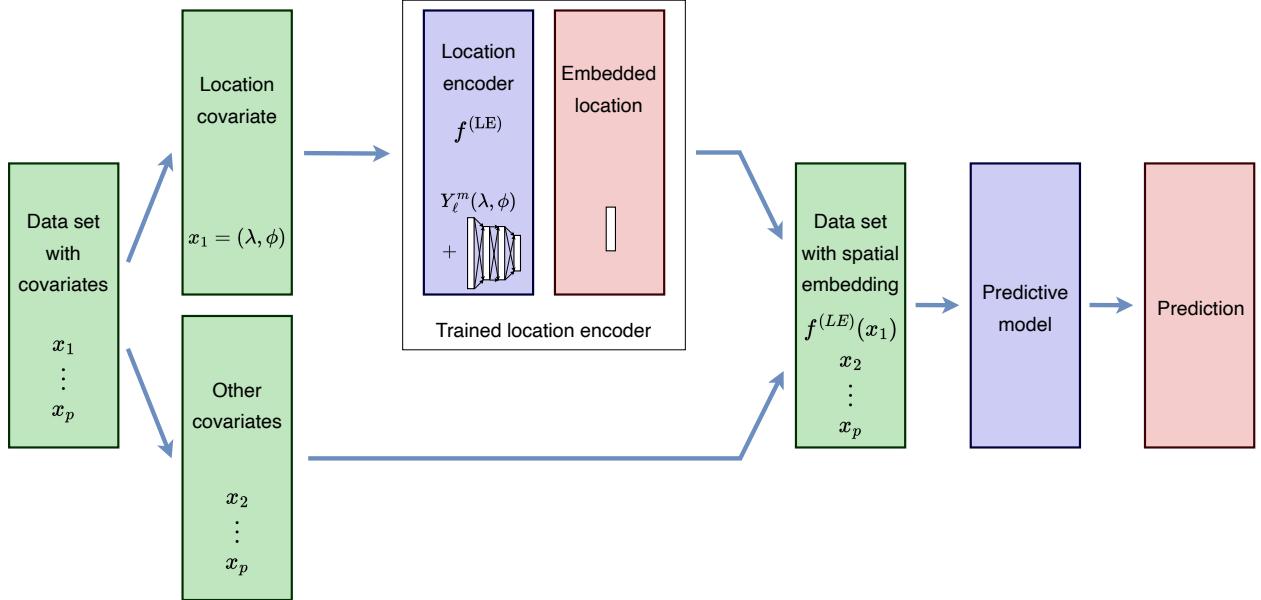


Figure 10: Application of the trained spatial embedding model in a case study where geographical coordinates are used as a covariate in a predictive model. By replacing the raw coordinate x_1 with its spatial embedding $f^{\text{LE}}(x_1)$, the model incorporates not only geographic position but also the spatial information observed at that location through the embedding.

3.4 Hyperparameter choices and trained model architectures

In the construction and training of our spatial embedding model architecture, several hyperparameters have to be specified. The OpenStreetMap encoder f^{OSM} consists of two convolutional layers, the first with 8 filters and the second with 16 filters, both with kernel size $k = 1$. The multi-view fusion encoder is implemented as a single dense layer consisting of 128 neurons. The

location encoder uses $L = 64$ spherical harmonic basis functions in the positional encoder, followed by a two-layer feed-forward network with 128 neurons each. Optimization is performed with the Adam algorithm (Kingma and Ba, 2015) using a learning rate of 10^{-5} and weight decay of 0.01. Training is carried out for 500 epochs with a batch size of 128. The contrastive loss employs a temperature parameter of $\tau = 0.07$. The hyperparameter choices either follow the values proposed in Klemmer et al. (2025) or were determined during preliminary testing in the model construction phase. More details can be found in Appendix B.

For the experimental study in Section 4, we study the effect of the overall embedding dimension d , and the intermediate dimensions d_{GS} and d_{OSM} of the Google Satellite and OpenStreetMap encoders on the resulting spatial embeddings. To this end, we construct five different embedding models, each corresponding to a different combination of d , d_{GS} , and d_{OSM} . Table 1 shows the values of d , d_{GS} , and d_{OSM} for each of the five embedding models, which are compared in the experimental study in Section 4.

Model name	Embedding dimension d	GS spatial view	OSM spatial view
EU8_GS32_OSM32	8	$d_{\text{GS}} = 32$	$d_{\text{OSM}} = 32$
EU16_GS32_OSM16	16	$d_{\text{GS}} = 32$	$d_{\text{OSM}} = 16$
EU16_OSM16	16	$d_{\text{GS}} = 0$	$d_{\text{OSM}} = 16$
EU32_GS96_OSM32	32	$d_{\text{GS}} = 96$	$d_{\text{OSM}} = 32$
EU64_GS64	64	$d_{\text{GS}} = 64$	$d_{\text{OSM}} = 0$

Table 1: Overview of the five spatial embedding models used in the experimental study. Each model varies the overall embedding dimension d and the intermediate dimensions d_{GS} and d_{OSM} , while all other hyperparameters are fixed. Model acronyms are constructed as follows: the prefix EU indicates that the model is trained on the European dataset introduced in Section 2; this is followed by the overall embedding dimension d ; then GS with the dimension of the Google Satellite encoder d_{GS} ; and finally OSM with the dimension of the OpenStreetMap encoder d_{OSM} .

All training is conducted using Python on a desktop computer with a single NVIDIA RTX 4070 GPU (12 GB), an Intel i9-13900K CPU, and 128GB of RAM. Training the model on the full European dataset required approximately 12 hours, depending on hyperparameters. Calculating the embeddings for a new coordinate, as illustrated in Figure 10, takes under 1 millisecond.

4 Experimental study

4.1 Real estate transactions dataset

We use a real estate transaction dataset published by Etalab (nd). Our dataset contains 200 000 property sales in France between 2017 and 2023. Each record in the dataset describes one transaction with the sale price and descriptive variables about the property. Table 2 provides an overview of the variables in the dataset, including their descriptions and whether they are continuous or categorical.

Variable	Original name (Etalab)	Description
price	valeur_fonciere	Transaction price at which the property was sold, as a continuous variable.
type	type_local	Type of property, as a categorical variable with three levels: house (type 1), apartment (type 2), or industrial/commercial (type 3).
land_use	code_nature_culture	Land-use classification code, as a categorical variable with 27 levels. See Appendix C for the codes and their meaning.
rooms	nombre_pieces_principale	Number of rooms in the property, as a continuous variable.
building_area	surface_reelle_bati	Surface area of the building in m^2 , as a continuous variable.
land_area	surface_terrain	Surface area of the land in m^2 , as a continuous variable.
date	date_mutation	Date of transaction, represented as three continuous variables: year, month, and day.
coordinates	latitude, longitude	Coordinates of the property, expressed as two continuous variables.

Table 2: Description of the variables in the French real estate transaction dataset. The column *Variable* lists the variable names used in this paper, while the column *Original name (Etalab)* shows the dataset field names as defined in Etalab ([nd](#)). Each description specifies whether the variable is continuous or categorical, and for categorical variables, the levels are listed.

The case study constructs a regression model using the variables in Table 2 as covariates with `price` as the response variable. The left-hand side of Figure 11 shows the distribution of the response variable in the dataset, and the right-hand side shows the average property value by postal code. All transaction prices lie between €100 000 and €2 000 000. The average transaction value of properties is €310 190, with 90% of properties having a transaction value between €112 000 and €800 000. The areas with a higher average transaction value include the Paris district, the south-east and south-west coast areas, and the French Alps bordering Switzerland. The centre of France features large areas with lower average transaction prices in the regions of Limousin, Auvergne, and Bourgogne, which are characterised by extensive areas of nature reserves and agricultural land. We observe a higher average transaction price in cities such as Orleans, Lyon, and Toulouse. The dataset does not contain properties in the departments Bas-Rhin, Haut-Rhin, and Moselle.

Figure 12 shows the distribution of the other variables in the French real estate dataset. The properties in the dataset are on average 135 m^2 , with 95% of properties being below 220 m^2 . The properties have on average four rooms, and a land area of 1012 m^2 . Around 81% of properties are classified as houses (Type 3), 12.5% are apartments (Type 1), and the others are industrial or commercial properties (Type 2). Most land use classifications are building lots (class S). More properties were sold in 2021 in comparison to 2020, and no sales were recorded in the latter half of 2022.

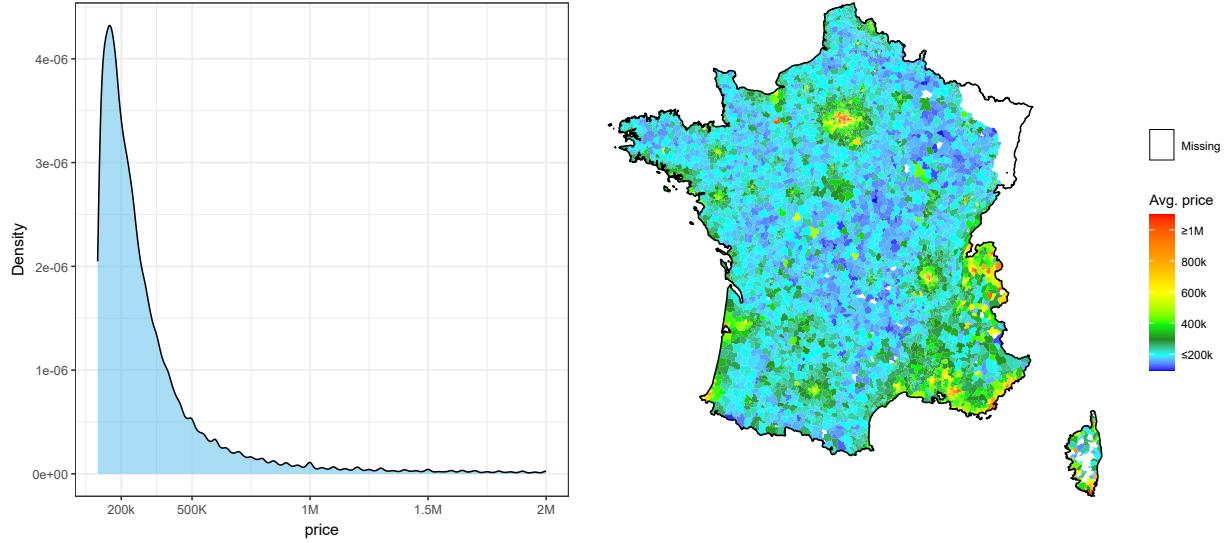


Figure 11: Plot of the property values in the dataset. The left-hand side shows a density plot of the property values. The right-hand side shows the average property value by postal code.

We denote the French real estate transaction dataset as $\mathcal{D}^{\text{FR}} = (\mathbf{x}_i, y_i)_{i=1}^{200\,000}$. Each index i represents one property in the dataset, with the vector $(x_{i,1}, \dots, x_{i,7})$ containing the value of each of the seven variables in the dataset, as described in Table 2, and y_i the value of the price of that property. Let the value $x_{i,7}$ be the latitude-longitude coordinates of the location of property i . When using a spatial embedding model, we replace the latitude-longitude coordinates of each property with the spatial embedding vector representing those coordinates, constructed by the location encoder f^{LE} from a trained spatial embedding model defined in Section 3.4.

4.2 Performative analysis of spatial embeddings

We fit a regression model on the French real estate dataset with the transaction price as the response variable and the variables in Table 2 as input features. We will compare the effect on predictive accuracy when using the latitude and longitude coordinates as input features versus replacing latitude-longitude with its spatial embedding vector. We construct both a generalised linear model (GLM), a generalised additive model (GAM), and a gradient boosting machine (GBM). The dataset is split into a randomly selected 80% of the observations as training data and the remaining 20% as an out-of-sample test set.

We fit the GLM using a Gaussian distribution and a logarithmic link function. The categorical input features are dummy encoded for use in the GLM. The levels with the largest amount of exposure are taken as reference level: `type1` and `land_useS`. The location variable is used as two continuous variables, `latitude` and `longitude`, and the date variable is used as three continuous variables, `year`, `month`, and `day`. Variable selection is done via a bidirectional stepwise variable selection process; see Hastie et al. (2001). The full model structure of the GLM using latitude and

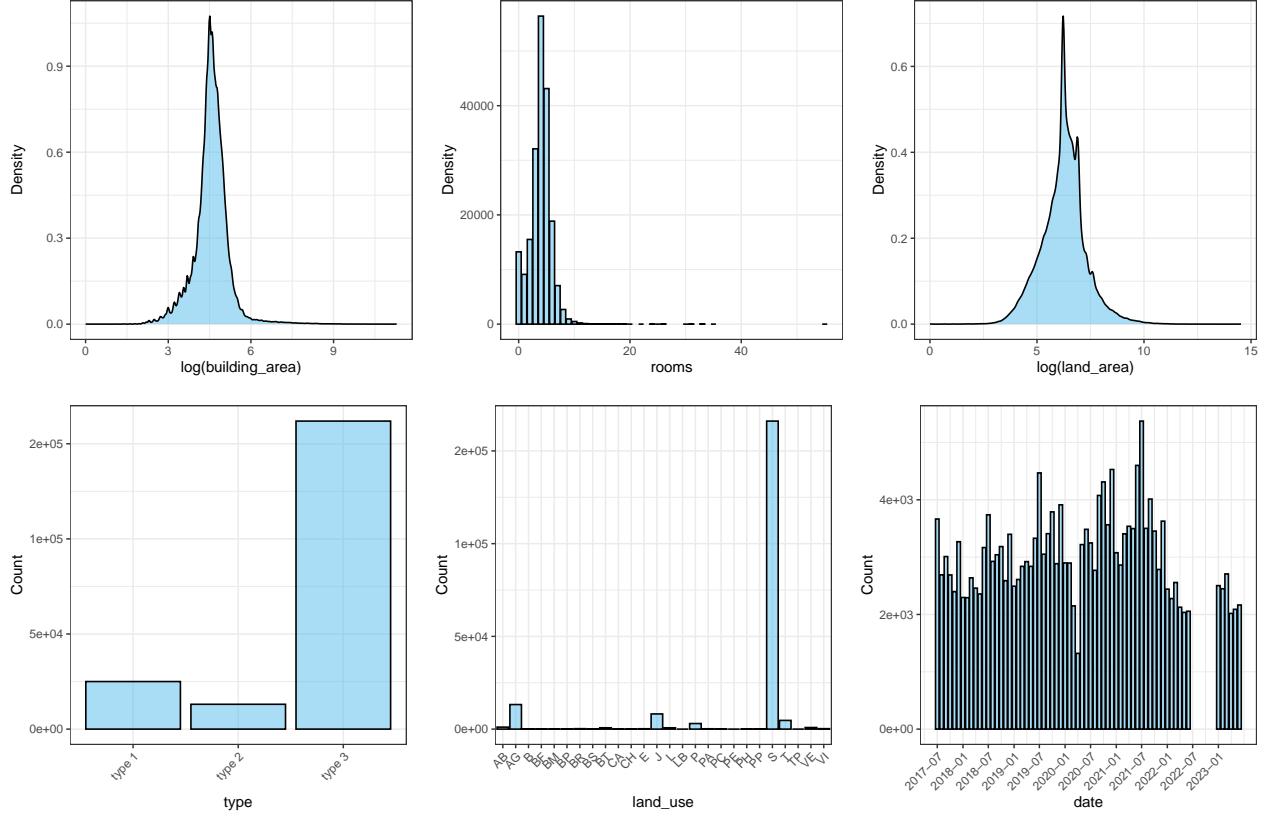


Figure 12: Exploratory data analysis of the French real-estate transaction dataset. From left to right, top row: density of $\log(\text{building_area})$, bar graph of rooms , and the density of $\log(\text{land_area})$. Bottom row: bar graph of type , bar graph of land_use , and a bar graph showing the number of transactions per month over time.

longitude as input variables is as follows:

$$\begin{aligned} \log(\mathbb{E}(\text{price})) = & \beta_0 + \beta_1 \text{type}_2 + \beta_2 \text{type}_3 + \beta_3 \text{land_use}_{\text{AB}} + \dots + \beta_{28} \text{land_use}_{\text{VI}} \\ & + \beta_{29} \text{rooms} + \beta_{30} \text{building_area} + \beta_{31} \text{land_area} \\ & + \beta_{32} \text{year} + \beta_{33} \text{month} + \beta_{34} \text{day} + \beta_{35} \text{latitude} + \beta_{36} \text{longitude}. \end{aligned} \quad (3)$$

For the GLM using the spatial embedding as input variable, the terms $\beta_{35} \text{latitude} + \beta_{36} \text{longitude}$ are replaced with $\beta_{35} \text{spatial_embedding}_1 + \dots + \beta_{34+d} \text{spatial_embedding}_d$, where $(\text{spatial_embedding}_1, \dots, \text{spatial_embedding}_d)$ is the spatial embedding vector of dimension d .

We use a Gaussian GAM with a logarithmic link function. The categorical variables are used similarly to those in the GLM, with dummy encoding and the same reference levels. For the continuous variables rooms , building_area , and land_area , univariate smooth functions f_1, f_2, f_3 are used. The variable date is treated as three continuous variables, year , month , and day , each with a univariate smooth function f_4, f_5, f_6 . The variable coordinates is used as two continuous variables, latitude and longitude , with a bivariate smooth function f_7 . The full model structure

of the GAM is as follows:

$$\begin{aligned} \log(\mathbb{E}(\text{price})) = & \beta_0 + \beta_1 \text{type}_2 + \beta_2 \text{type}_3 + \beta_3 \text{land_use}_{AB} + \dots + \beta_{28} \text{land_use}_{VI} \\ & + f_1(\text{rooms}) + f_2(\text{building_area}) + f_3(\text{land_area}) \\ & + f_4(\text{year}) + f_5(\text{month}) + f_6(\text{day}) + f_7(\text{latitude}, \text{longitude}). \end{aligned} \quad (4)$$

For the GAM using the spatial embedding as input variable, the term $f_7(\text{latitude}, \text{longitude})$ is replaced by $f_7(\text{spatial_embedding}_1) + \dots + f_{6+d}(\text{spatial_embedding}_d)$.

For the GBM, we use the XGBoost implementation (Chen and Guestrin, 2016). The categorical input features are one-hot encoded for use in the GBM. The optimal tuning parameter values are determined via five-fold cross-validation. We tune the number of trees from the values $\{100, 500, 1000, 5000\}$, the depth of each tree from the values $\{3, 5, 7\}$, and the learning rate from the values $\{0.01, 0.001\}$. We assess model performance for the GLM, GAM, and GBM using the mean squared error (MSE) of log-transformed house prices on the test set. The same test set is used for all three models to ensure a fair comparison.

Table 3 shows the out-of-sample losses for the GLM, GAM, and GBM. The first column specifies how the location of the property was used in the model: either using latitude-longitude as input features (no spatial embedding), or a spatial embedding vector as an input feature, in which case the name of the embedding model follows the acronyms given in Table 1. We indicate the lowest out-of-sample loss for each column in bold. We observe that for the GLM, GAM, and GBM, using spatial embeddings instead of latitude-longitude coordinates as input features results in improved out-of-sample loss. The larger the dimension of the embeddings, the lower the out-of-sample loss, which is to be expected as larger embedding dimensions allow for more spatial information to be stored in the embeddings.

Spatial treatment	GLM	GAM	GBM
No spatial embedding	0.3122	0.2362	0.1415
EU8_GS32_OSM32	0.3070	0.2147	0.1344
EU16_GS32_OSM16	0.2861	0.1978	0.1332
EU16_OSM16	0.2775	0.1929	0.1328
EU32_GS96_OSM32	0.2646	0.1829	0.1325
EU64_GS64	0.2473	0.1762	0.1310

Table 3: Comparison between models using latitude-longitude and spatial embedding as input features. Results show out-of-sample losses, measured in mean-square error on the log transform of the transaction price. For spatial treatment, the model acronyms are used as defined in Table 1. The lowest out-of-sample loss per column is indicated in bold.

4.3 Model interpretation of spatial effect

Variable importance We use the permutational variable importance tool from Olden et al. (2004) to compare the importance of each variable in our dataset on the prediction of real estate price. The variable importance is measured as the average change in prediction when randomly

permuting the values of an input feature. For a prediction model $f_{\text{pred}}(\cdot)$, and an input feature x_j in the dataset \mathcal{D}^{FR} , the importance of that feature is defined as

$$\text{VIP}_{x_j} = \frac{1}{|\mathcal{D}^{\text{FR}}|} \sum_{\mathbf{x}_i \in \mathcal{D}^{\text{FR}}} \text{abs}(f_{\text{pred}}(x_{i,1}, \dots, x_{i,j}, \dots, x_{i,7}) - f_{\text{pred}}(x_{i,1}, \dots, \tilde{x}_{i,j}, \dots, x_{i,7})), \quad (5)$$

where $\tilde{x}_{i,j}$ is a random permutation of all values of that input variable. When calculating the variable importance of the property’s location, the latitude–longitude pairs are permuted together rather than the latitude and longitude values separately. Similarly, when using the spatial embedding to represent location, the entire embedding vectors are permuted, not the individual elements within each vector. We calculate the variable importance for each of the input variables in Table 2. Figure 13 shows the variable importance for the GLM on the left, the GAM in the middle, and the GBM on the right. Dark blue shows the effect of the model with spatial embeddings, whereas light blue shows the effect of the model using latitude-longitude coordinates. For the spatial embedding, we utilised the location encoder from the trained EU32_GS96_OSM32 model, as described in Section 3.4. We observe in the GLM that the coordinates of the property are the third most important variable. This is because the GLM can only fit two coefficients, one for latitude and one for longitude, which is not sufficient to capture any nuanced spatial effects. Using spatial embeddings, however, the embedding vector is seen as the most important variable in the model. For both the GAM and the GBM, location is the most important variable, whether used via latitude and longitude coordinates or spatial embeddings. For the GLM and the GAM, the effect of permuting the spatial embeddings is so significant that the relative importance of other variables is no longer visible in Figure 13. This is not the case for the GBM model, where the importance of non-location variables stays relatively equal between the model using latitude and longitude versus the model using spatial embeddings.

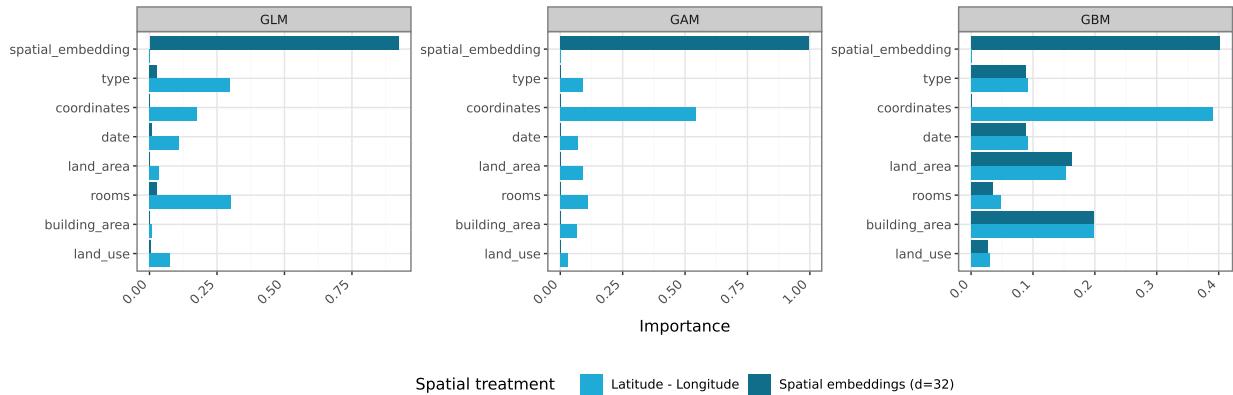


Figure 13: Variable importance as measured in the GLMs (left side), the GAM (middle), and the GBM models (right side). In dark blue, we show the importance of each variable as measured in the model with spatial embeddings. In light blue, we show the importance measured in the model using latitude-longitude coordinates.

Partial dependence We use partial dependence plots (Hastie et al., 2009) to examine the relationship between the predicted real estate price and the property’s location in the pricing model.

For a prediction model $f_{\text{pred}}(\cdot)$, where input feature $x_{i,7}$ represents the latitude-longitude coordinates, we calculate the partial dependence effect of a location (λ, ϕ) as

$$\text{PD}(\lambda, \phi) = \frac{1}{|\mathcal{D}^{\text{FR}}|} \sum_{\mathbf{x}_i \in \mathcal{D}^{\text{FR}}} f_{\text{pred}}(x_{i,1}, \dots, x_{i,6}, (\lambda, \phi)). \quad (6)$$

When using a prediction model with the spatial embedding of latitude-longitude coordinates as input feature, as calculated via the encoder f^{LE} , the partial dependence effect of the location (λ, ϕ) is calculated as

$$\text{PD}(\lambda, \phi) = \frac{1}{|\mathcal{D}^{\text{FR}}|} \sum_{\mathbf{x}_i \in \mathcal{D}^{\text{FR}}} f_{\text{pred}}(x_{i,1}, \dots, x_{i,6}, f^{\text{LE}}(\lambda, \phi)). \quad (7)$$

Let $\{(\lambda_c, \phi_c), c = 1, \dots, 6048\}$ denote the latitude-longitude coordinates of the centroids of the 6048 postal codes in France. For each coordinate pair, we compute the partial dependence effect, which results in a spatial map showing how location influences the predicted real estate price.

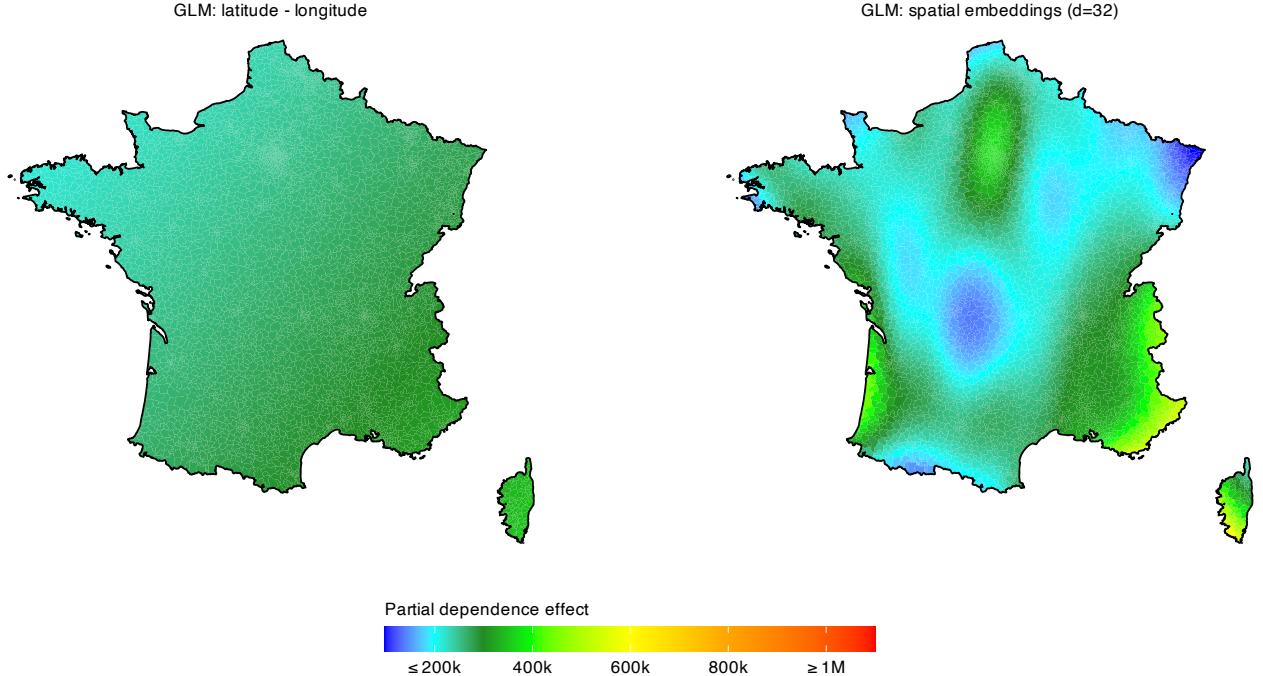


Figure 14: Partial dependence effect of location on real estate price, calculated per postal code for predictions from the GLMs. The left-hand side is the model using latitude-longitude coordinates as input features, and the right-hand side is the model using spatial embeddings instead of latitude-longitude. We use the EU32_GS96_OSM32 model to extract spatial embeddings.

Figure 14 shows the partial dependence effect as calculated with the GLM. The left-hand side shows the effect of the GLM using latitude and longitude as input features as defined in Equation (3); the right-hand side shows the GLM using the spatial embedding vector as input features. The GLM using latitude and longitude shows a negligible partial dependence effect, with each postal code having a very similar average predicted price. This is in line with the variable importance

plots in Figure 13, which indicate that the variables latitude and longitude are not important in this GLM. By contrast, the partial dependence effect for the GLM using spatial embeddings shows higher predicted prices around Paris and in the south-east and south-west regions of France, with a further effect visible on the island of Corsica. This demonstrates that the GLM can capture more nuanced patterns through spatial embeddings than it can solely based on latitude and longitude coordinates.

Figure 15 shows the partial dependence effect per postal code as calculated with the GAM. The left-hand side displays the GAM using latitude and longitude as input, as described in Equation (4). Compared to the GLM, the bivariate smoothing on the coordinates produces a more nuanced spatial pattern, with higher predicted prices around Paris and in the south-east and south-west regions. The right-hand side shows the GAM using spatial embeddings as input. In this case, the model captures more complex spatial variation, distinguishing between urban centres and rural areas.

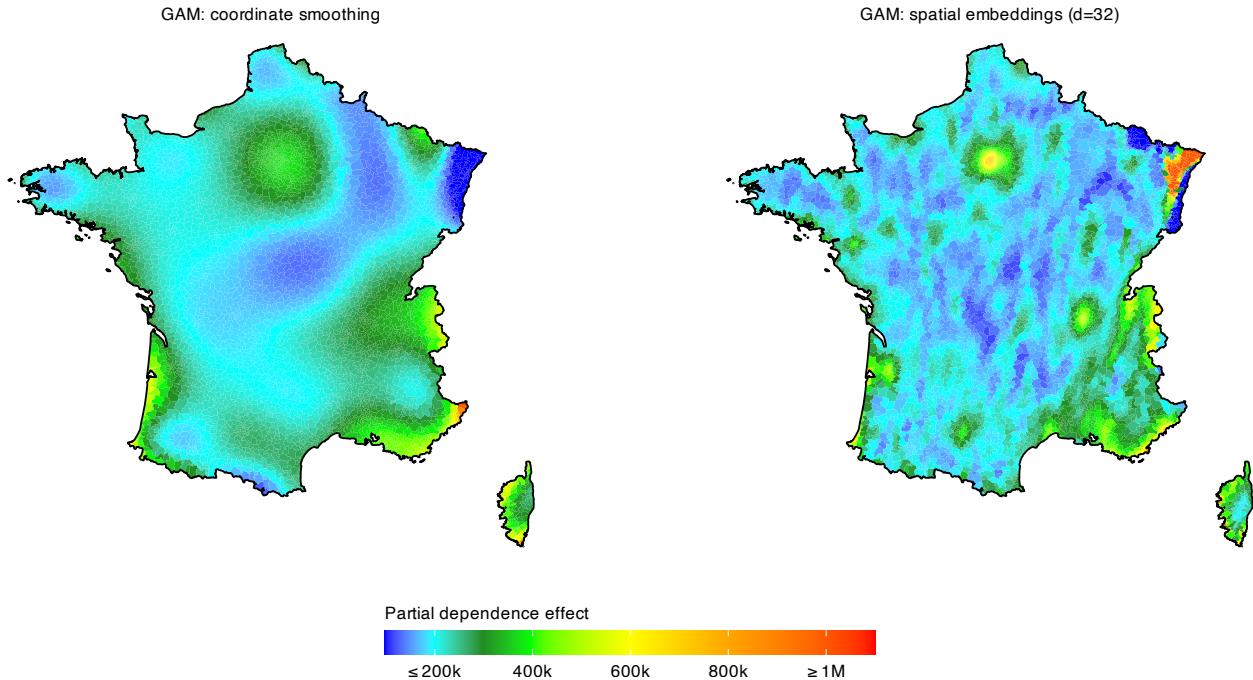


Figure 15: Partial dependence effect of location on real estate price, calculated per postal code for predictions from the GAM. The left-hand side is the model using latitude-longitude coordinates as input features, and the right-hand side is the model using spatial embeddings instead of latitude-longitude. We use the EU32_GS96_OSM32 model to extract spatial embeddings.

Figure 16 shows the partial dependence effect of the GBM using latitude-longitude as input features on the left and using the spatial embeddings on the right-hand side. For both the GBM using latitude-longitude and the GBM using spatial embeddings, we see a difference between cities in France and rural areas. The GBM using latitude-longitude, however, shows a square pattern caused by the linear splits made by the regression trees in the GBM (Henckaerts et al., 2021; Geerts et al., 2024). The partial dependence effect from the GBM using spatial embeddings does not show these square artifacts and shows the shapes of cities more clearly.

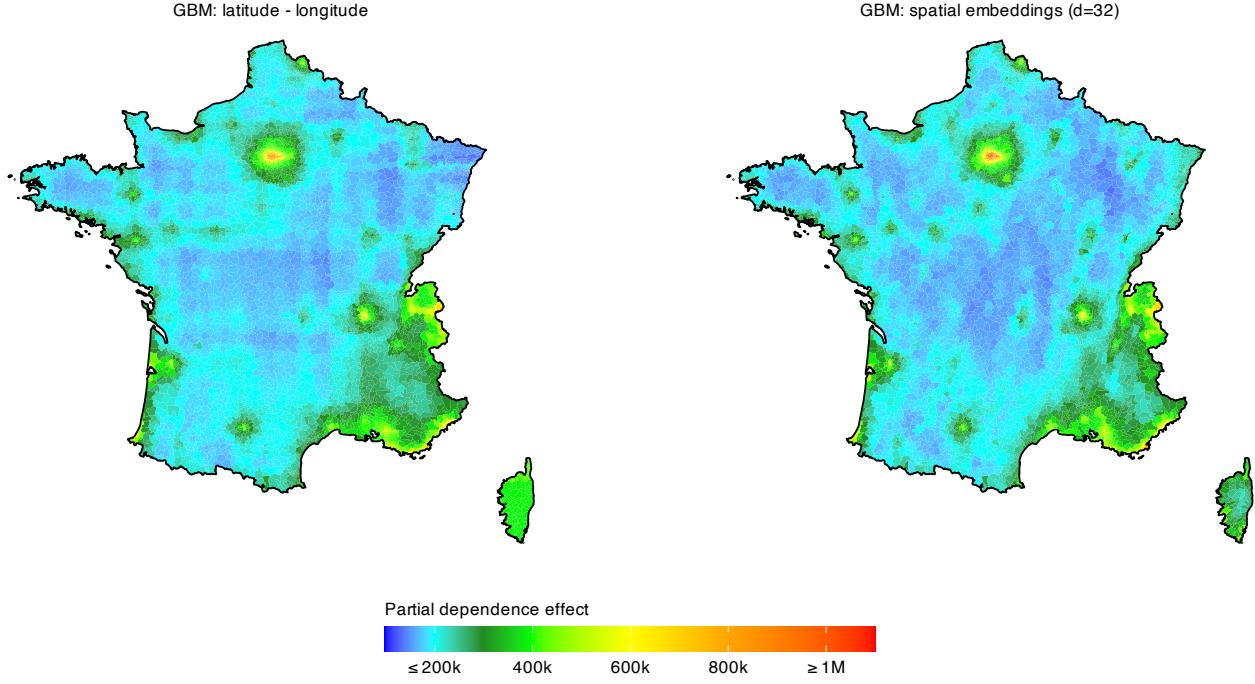


Figure 16: Effect of location on average predicted real estate price, calculated per postal code for predictions from the GBMs. The left-hand side is the model using latitude-longitude coordinates as input features, and the right-hand side is the model using spatial embeddings instead of latitude-longitude. We use the EU32_GS96_OSM32 model to extract spatial embeddings.

Fine-grained partial dependence effects Instead of one coordinate pair per postal code, we define a fine-grained set of locations in the area between Lyon and Switzerland. We randomly sample a set of 10 000 locations in this area. By calculating the partial dependence effect for this set of locations, we identify fine-grained effects in the relationship between location and predicted prices. We calculate the partial dependence effect for each location following Equation (7) using the GBM.

Figure 17 shows the partial dependence effects. The left-hand side displays results at the postal code level, while the right-hand side shows the fine-grained set of locations. Three localised patterns appear in the fine-grained map that are not visible at the postal code level, highlighted by red ovals. The first is a small, oblong area with higher predicted prices, corresponding to the touristically popular northwestern slopes of the Aravis range. The second is a wide oblong area with lower predicted prices, coinciding with the La Lauzière range. The La Lauzière mountain range is primarily not used for tourism and houses only a few resorts or ski areas. The third is an angular pattern of higher predicted prices matching the Le Trois Vallées area. This is the largest connected ski area in the world and is highly touristic. These examples demonstrate that the GBM has fitted a relationship between real estate prices and the spatial information captured in the spatial embeddings.

Transfer learning based on spatial embeddings We compare the GBM's capacity for spatial transfer learning between the model using latitude-longitude coordinates and the model using spatial embeddings. Spatial transfer learning refers to a model's ability to produce accurate pre-

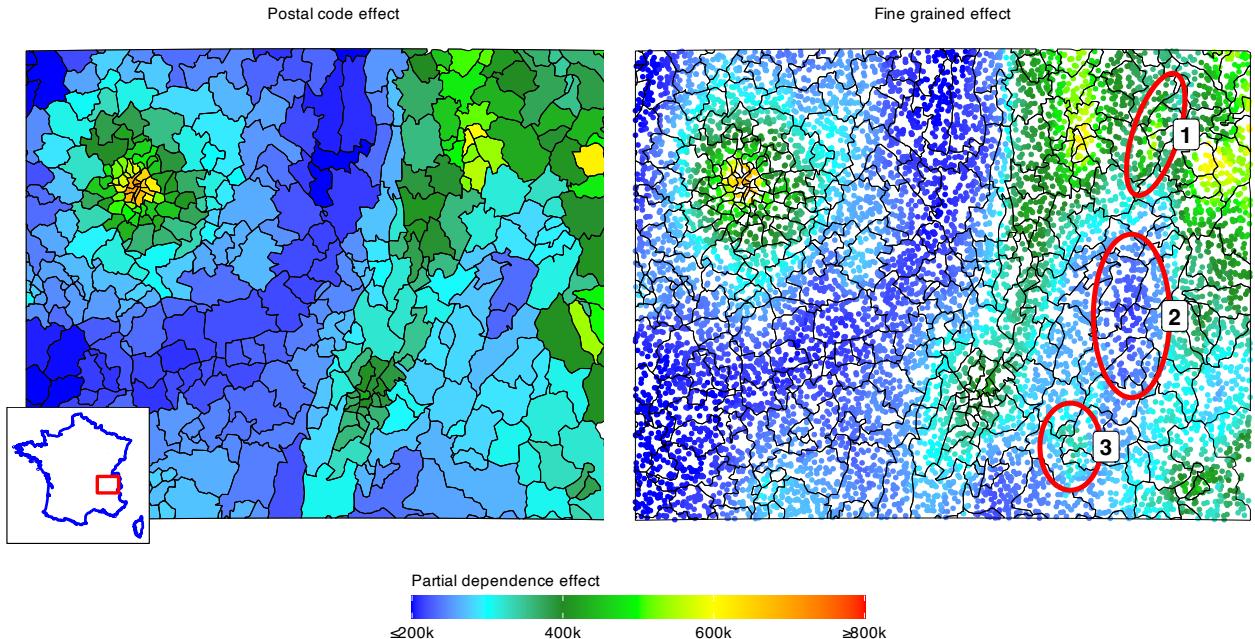


Figure 17: Partial dependence effect from the GBM model using spatial embeddings calculated in a selected area of France between Lyon and Switzerland. The used spatial embedding model is the EU32_GS96_OSM32 model. The selected area is indicated on the small map of France on the bottom left. The left-hand side shows the partial dependence effect calculated on the level of postal codes, and the right-hand side shows the partial dependence effect calculated on 10 000 randomly sampled points in the selected area.

dictions in regions outside the spatial domain of its training data. The coordinate-based model can only extrapolate based on geographic position, whereas the embedding-based model can adapt more effectively by leveraging the spatial characteristics represented by the embeddings in the unseen region. We examine this effect by analysing the partial dependence for postal codes in the departments Bas-Rhin, Haut-Rhin, and Moselle, which are excluded from the French real estate dataset.

The first panel of Figure 18 shows the empirical mean transaction price per postal code in the region between Nancy and the German/Swiss border; a map of France in the lower left highlights the selected area. A grey line marks the spatial extent of the training data, and postal codes in white indicate areas absent from the dataset. The second panel shows the partial dependence (PD) for this region obtained from the GBM using latitude and longitude as inputs. Because the excluded departments lie beyond the range of latitude and longitude values present in the training data, the GBM cannot form new splits in these areas. Instead, the model extrapolates by assigning these unseen locations to the nearest regions learnt during training, resulting in an unrealistic, piecewise-constant partial dependence effect. The third panel shows the partial dependence effect obtained from the GBM using spatial embeddings as input features. Here, the model captures a lower partial dependence effect in the central part of the region and a higher effect along the Swiss border. Both patterns are intuitive: the central area is predominantly rural with lower population density, whereas larger cities, such as Strasbourg, are located along the border. By using spatial embeddings, the model associates predictions with spatial characteristics rather than with absolute

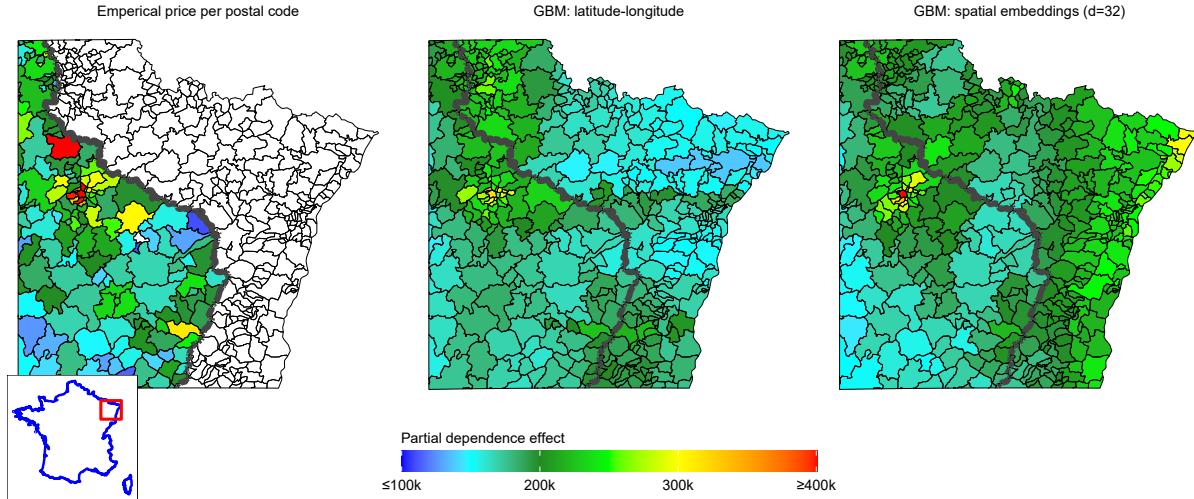


Figure 18: Selected area in France between the city of Nancy and the German border, as indicated on the inset map (bottom left). The left panel shows the empirically observed average transaction price per postal code, where white corresponds to postal codes without observations in the dataset, namely in the departments Bas-Rhin, Haut-Rhin, and Moselle. The middle panel shows the partial dependence effect calculated with the GBM using latitude and longitude as input. The right panel shows the partial dependence effect from the GBM using spatial embeddings. The grey line marks the boundary between regions included in the training data and regions without observations.

coordinates. This representation supports spatial transfer learning, enabling the model to generalise to unseen regions by exploiting similarities in spatial structure through the spatial embeddings.

5 Conclusion

In this paper, we introduced a framework for constructing spatial embeddings using contrastive learning in a multi-view setting. To demonstrate the multi-view spatial embedding model, we combined satellite imagery and OpenStreetMap tags to learn low-dimensional representations of geographic locations. These embeddings are trained to align with coordinate-based encodings, allowing any dataset with latitude–longitude information to be enriched with spatial context without requiring direct access to the original spatial inputs. The paper is accompanied by a GitHub repository at <https://github.com/freekholvoet/MultiviewSpatialEmbeddings>, which provides the code for training the spatial embedding model. The trained models defined in Table 1 are available in Holvoet (2025).

We demonstrated the practical value of these embeddings in a real estate pricing application. Replacing raw coordinates with spatial embeddings consistently improved predictive performance across GLMs, GAMs, and GBMs. The embeddings captured spatial structure more effectively than traditional location features and provided interpretable spatial effects. Variable importance and

partial dependence analyses showed that the embeddings learnt meaningful geographic patterns, including urban–rural differences, regional economic variation, and fine-grained effects related to topography and land use. These effects are difficult to model using raw coordinates or categorical location indicators such as postal codes. Moreover, the use of spatial embeddings enabled models to generalise to regions not seen during training by associating predictions with spatial characteristics rather than coordinates. This property is particularly relevant for insurance applications, where actuaries often need to assess risk in new or underrepresented areas.

From an actuarial modelling perspective, spatial embeddings offer a flexible and scalable way to incorporate geographic information into pricing models. They overcome the limitations of entity embeddings and spatial smoothing methods and can be used in combination with other modelling techniques. The embeddings are smooth, transferable, and computationally efficient, making them suitable for large-scale applications in insurance analytics. A further advantage of the proposed framework is that it can be trained on modest computational hardware and easily adapted to new regions or additional spatial data sources, making it practical for in-house deployment at insurance companies and other organisations.

This work represents a step toward integrating the principles of foundation models into actuarial science. By demonstrating how multimodal spatial pretraining produces reusable embeddings that enhance predictive accuracy and generalise across space, we show that actuaries can benefit from adopting foundation-model approaches. Just as language and vision models have been transferred to countless downstream tasks, spatial foundation models can serve as general-purpose representations that can be adapted for pricing, reserving, portfolio analytics, and emerging climate-related risks. While recent global-scale spatial foundation models, such as AlphaEarth, illustrate what is possible at the planetary scale, our work highlights a complementary direction: lightweight, region-specific spatial embeddings that can be retrained quickly, incorporate organisation-specific spatial views, and capture fine-grained local structure that global models may overlook. We encourage the actuarial community to explore and experiment with such models, to evaluate how pretrained spatial representations can be transferred across products, regions, and lines of business, and to impute domain expertise in shaping foundation models that are both scientifically rigorous and practically relevant for risk management.

Acknowledgements

Katrien Antonio gratefully acknowledges funding from the FWO and Fonds De La Recherche Scientifique - FNRS (F.R.S.-FNRS) under the Excellence of Science (EOS) program, project ASTeRISK Research Foundation Flanders [grant number 40007517]. The authors gratefully acknowledge support from the FWO network W001021N. Christopher Blier-Wong acknowledges financial support from the Natural Sciences and Engineering Research Council of Canada (RGPIN-2025-06879).

Declaration of interest statement

The authors declare no potential conflicts of interest.

Supplemental material

The code to train the model architecture can be found via the Github repository: <https://github.com/freekholvoet/MultiviewSpatialEmbeddings>. The trained models, as defined in Table 1, are available in Holvoet (2025).

References

- Avanzi, B., Taylor, G., Wang, M., and Wong, B. (2024). Machine learning with high-cardinality categorical features in actuarial applications. *ASTIN Bulletin : The Journal of the IAA*, 54(2):213–238.
- Blier-Wong, C., Baillargeon, J.-T., Cossette, H., Lamontagne, L., and Marceau, E. (2020). Encoding neighbor information into geographical embeddings using convolutional neural networks. In *Proceedings of the Thirty-Third International FLAIRS Conference*. <https://aaai.org/papers/15-flairs-2020-18400/>.
- Blier-Wong, C., Baillargeon, J.-T., Cossette, H., Lamontagne, L., and Marceau, E. (2021). Rethinking representations in P&C actuarial science with deep neural networks. *arXiv:2102.05784*. <https://arxiv.org/abs/2102.05784>.
- Blier-Wong, C., Cossette, H., Lamontagne, L., and Marceau, E. (2022). Geographic ratemaking with spatial embeddings. *ASTIN Bulletin: The Journal of the IAA*, 52(1):1–31.
- Boskov, M. and Verrall, R. (1994). Premium rating by geographic area using spatial models. *ASTIN Bulletin: The Journal of the IAA*, 24(1):131–143.
- Brown, C. F., Kazmierski, M. R., Pasquarella, V. J., Ruckridge, W. J., Samsikova, M., Zhang, C., Shelhamer, E., Lahera, E., Wiles, O., Ilyushchenko, S., Gorelick, N., Zhang, L. L., Alj, S., Schechter, E., Askay, S., Guinan, O., Moore, R., Boukouvalas, A., and Kohli, P. (2025). AlphaEarth foundations: An embedding field model for accurate and efficient global mapping from sparse label data. *arXiv:2507.22291*. <https://arxiv.org/abs/2507.22291>.
- Chen, T. and Guestrin, C. (2016). XGBoost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 785–794. <https://dl.acm.org/doi/10.1145/2939672.2939785>.
- Delong, L. and Kozak, A. (2023). The use of autoencoders for training neural networks with mixed categorical and numerical features. *ASTIN Bulletin: The Journal of the IAA*, 53(2):213–232.
- Dimakos, X. K. and Di Rattalma, A. F. (2002). Bayesian premium rating with latent structure. *Scandinavian Actuarial Journal*, 2002(3):162–184.
- Etalab (n.d.). Demandes de valeurs foncières géolocalisées. <https://www.data.gouv.fr/en/datasets/demandes-de-valeurs-foncieres-geolocalisees/>. [Data set].
- Eurostat (2021). Local administrative units (lau). <https://ec.europa.eu/eurostat/web/nuts/local-administrative-units>. Accessed: 2024-08-28.
- Gao, L. and Shi, P. (2022). Leveraging high-resolution weather information to predict hail damage claims: A spatial point process for replicated point patterns. *Insurance: Mathematics and Economics*, 107:161–179.
- Geerts, M., vanden Broucke, S., and De Weerdt, J. (2024). GeoRF: a geospatial random forest. *Data Mining and Knowledge Discovery*, 38(6):3414–3448.

- Google (n.d.). Google Maps Static API. <https://developers.google.com/maps/documentation/maps-static>. Accessed: 2024-10-28.
- Gschlößl, S. and Czado, C. (2007). Spatial modelling of claim frequency and claim size in non-life insurance. *Scandinavian Actuarial Journal*, 2007(3):202–225.
- Guo, C. and Berkhahn, F. (2016). Entity embeddings of categorical variables. *arXiv preprint arXiv:1604.06737*. <https://arxiv.org/abs/1604.06737>.
- Hastie, T., Tibshirani, R., and Friedman, J. H. (2001). *The elements of statistical learning*. Springer Series in Statistics. Springer New York, New York, NY.
- Hastie, T., Tibshirani, R., and Friedman, J. H. (2009). *The elements of statistical learning: data mining, inference, and prediction*, volume 2 of *Springer Series in Statistics*. Springer New York, New York, NY.
- He, K., Fan, H., Wu, Y., Xie, S., and Girshick, R. (2020). Momentum contrast for unsupervised visual representation learning. In *Proceedings of the 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9726–9735. <https://ieeexplore.ieee.org/document/9157636>.
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778. <https://ieeexplore.ieee.org/document/7780459>.
- Henckaerts, R., Antonio, K., Clijsters, M., and Verbelen, R. (2018). A data driven binning strategy for the construction of insurance tariff classes. *Scandinavian Actuarial Journal*, 2018(8):681–705.
- Henckaerts, R., Côté, M.-P., Antonio, K., and Verbelen, R. (2021). Boosting insights in insurance tariff plans with tree-based machine learning methods. *North American Actuarial Journal*, 25(2):255–285.
- Holvoet, F. (2025). Spatial embeddings via multiview contrastive learning. [Pretrained spatial embedding models], https://huggingface.co/FreekH/multiview_spatial_embedding, DOI 10.57967/hf/7009.
- Holvoet, F., Antonio, K., and Henckaerts, R. (2025). Neural networks for insurance pricing with frequency and severity data: A benchmark study from data preprocessing to technical tariff. *North American Actuarial Journal*, 29(3):519–562.
- Kingma, D. P. and Ba, J. L. (2015). Adam: A method for stochastic optimization. In *Proceedings of the International Conference on Learning Representations (ICLR)*. <https://dblp.org/rec/journals/corr/KingmaB14.html>.
- Klemmer, K., Rolf, E., Robinson, C., Mackey, L., and Rußwurm, M. (2025). SatCLIP: Global, General-Purpose Location Embeddings with Satellite Imagery. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39(4), pages 4347–4355. <https://ojs.aaai.org/index.php/AAAI/article/view/32457>.
- Klemmer, K., ROLF, E., Robinson, C., Mackey, L., and Rußwurm, M. (2023). Towards global, general-purpose pretrained geographic location encoders. In *Proceedings of the NeurIPS 2023 Workshop on Tackling Climate Change with Machine Learning*. <https://www.climatechange.ai/papers/neurips2023/114>.
- Mai, G., Huang, W., Sun, J., Song, S., Mishra, D., Liu, N., Gao, S., Liu, T., Cong, G., Hu, Y., Cundy, C., Li, Z., Zhu, R., and Lao, N. (2024). On the opportunities and challenges of foundation models for GeoAI (vision paper). *ACM Transactions on Spatial Algorithms and Systems*, 10(2):1–46.

- Mai, G., Janowicz, K., Hu, Y., Gao, S., Yan, B., Zhu, R., Cai, L., and Lao, N. (2022). A review of location encoding for GeoAI: Methods and applications. *International Journal of Geographical Information Science*, 36(4):639–673.
- Mai, G., Janowicz, K., Yan, B., Zhu, R., Cai, L., and Lao, N. (2020). Multi-scale representation learning for spatial feature distributions using grid cells. In *Proceedings of the International Conference on Learning Representations*. https://iclr.cc/virtual_2020/poster_rJljd4KDH.html.
- Mai, G., Xie, Y., Jia, X., Lao, N., Rao, J., Zhu, Q., Liu, Z., Chiang, Y.-Y., and Jiao, J. (2025). Towards the next generation of geospatial artificial intelligence. *International Journal of Applied Earth Observation and Geoinformation*, 136: Article 104368.
- Moussaid, S. (2023). *Convolutional Regional Autoencoders for Spatial Embedding*. Master's thesis, KU Leuven, Leuven.
- Olden, J. D., Joy, M. K., and Death, R. G. (2004). An accurate comparison of methods for quantifying variable importance in artificial neural networks using simulated data. *Ecological Modelling*, 178(3-4):389–397.
- OpenStreetMap (n.d.). Openstreetmap. <https://www.openstreetmap.org>. Accessed: 2024-08-30.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., and Sutskever, I. (2021). Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139, pages 8748–8763. <https://proceedings.mlr.press/v139/radford21a.html>.
- Richman, R. (2021). AI in actuarial science – a review of recent advances – part 1. *Annals of Actuarial Science*, 15(2):207–229.
- Richman, R. and Wüthrich, M. V. (2024). High-cardinality categorical covariates in network regressions. *Japanese Journal of Statistics and Data Science*, 7:921–965.
- Robben, J., Antonio, K., and Kleinow, T. (2025). The short-term association between environmental variables and mortality: evidence from europe. *Journal of the Royal Statistical Society Series A: Statistics in Society*, page qnaf052.
- Rußwurm, M., Klemmer, K., Rolf, E., Zbinden, R., and Tuia, D. (2024). Geographic location encoding with spherical harmonics and sinusoidal representation networks. In *Proceedings of the International Conference on Learning Representations*. <https://iclr.cc/virtual/2024/poster/18690>.
- Shi, P. and Shi, K. (2017). Territorial risk classification using spatially dependent frequency-severity models. *ASTIN Bulletin: The Journal of the IAA*, 47(2):437–465.
- Shi, P. and Shi, K. (2023). Non-life insurance risk classification using categorical embedding. *North American Actuarial Journal*, 27(3):579–601.
- Sitzmann, V., Martel, J., Bergman, A., Lindell, D., and Wetzstein, G. (2020). Implicit neural representations with periodic activation functions. *Advances in Neural Information Processing Systems*, 33:7462–7473.
- Snyder, W. E., Qi, H., and Sander, W. A. (1999). Coordinate system for hexagonal pixels. In *Proceedings of Medical Imaging 1999: Image Processing*, volume 3661, pages 716 – 727. <https://www.spiedigitallibrary.org/conference-proceedings-of-spie/3661/0000/Coordinate-system-for-hexagonal-pixels/10.1117/12.348629.short>.
- Tarasiou, M., Chavez, E., and Zafeiriou, S. (2023). ViTs for SITS: Vision Transformers for Satellite

- Image Time Series. In *Proceedings of the 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10418–10428. IEEE. <https://ieeexplore.ieee.org/document/10204199>.
- Taylor, G. (2001). Geographic premium rating by Whittaker spatial smoothing. *ASTIN Bulletin: The Journal of the IAA*, 31(1):147–160.
- Taylor, G. C. (1989). Use of spline functions for premium rating by geographic area. *ASTIN Bulletin: The Journal of the IAA*, 19(1):91–122.
- Uber Technologies Inc. (n.d.). H3: A hexagonal hierarchical spatial index. <https://h3geo.org>. Accessed: 2024-08-30.
- Wang, Y., Braham, N. A. A., Xiong, Z., Liu, C., Albrecht, C. M., and Zhu, X. X. (2023). SSL4EO-S12: A large-scale multimodal, multitemporal dataset for self-supervised learning in earth observation [software and data sets]. *IEEE Geoscience and Remote Sensing Magazine*, 11(3):98–106.
- Wilsens, P., Antonio, K., and Claeskens, G. (2024). Reducing the dimensionality and granularity in hierarchical categorical variables. *Advances in Data Analysis and Classification*, pages 1–32.
- Yin, Y., Liu, Z., Zhang, Y., Wang, S., Shah, R. R., and Zimmermann, R. (2019). GPS2Vec: Towards generating worldwide GPS embeddings. In *Proceedings of the 27th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems*, pages 416–419. <https://dl.acm.org/doi/10.1145/3347146.3359067>.

A OpenStreetMap tags per category

Table 4 lists the OpenStreetMap tags used in each of the six categories as defined in Section 2. Each category lists the tag type, as defined by OpenStreetMap, along with the names of the tags collected per type.

Category	OpenStreetMap tag type	Tag names
shop	amenity	shops
tourism	amenity	tourism
food&drink	amenity	bar, cafe, pub, restaurant, fast food, food court, events venue, nightclub
health	amenity	clinic, dentist, doctors, hospital, pharmacy, social facility
	healthcare	clinic, doctor, hospital, pharmacy
education	amenity	college, kindergarten, research institute, school, university
public services	amenity	fire station, police

Table 4: The OpenStreetMap tags, with the tag type, for each of the six categories defined in Section 2.

B Hyperparameter selection

Hyperparameter selection for the spatial embedding model was carried out during model construction. As the embeddings are trained in a self-supervised manner, their quality cannot be directly evaluated using standard predictive metrics. Instead, the aim was to identify configurations that produced stable training behaviour and embeddings containing spatially informative structure.

A range of model configurations was tested for the OpenStreetMap encoder, location encoder, and fusion layer, varying the number of filters, neurons, and learning rates. Each configuration was pretrained for a fixed number of epochs and evaluated on two criteria. First, the contrastive loss (CLIP loss) was monitored to ensure stable convergence. Second, the resulting embeddings were applied in small-scale downstream tasks, such as logistic regression for country classification and a regression model for population density prediction. For both downstream tasks, we compared the performance of models with the embedding vector as input versus those with latitude and longitude coordinates. Configurations showing both stable loss behaviour and high downstream classification accuracy were retained.

The final hyperparameter values, reported in Section 3.4, combine established settings from Klemmer et al. (2025) with those identified through these preliminary tests. Smaller convolutional architectures (8 and 16 filters of size $k = 1$) and a 128-neuron fusion layer yielded consistent results, while larger models offered limited improvement and less stable optimisation. As in most self-supervised learning setups, the hyperparameter selection is therefore empirical, guided by training stability and representational quality rather than formal optimisation.

C Land use classification codes

Table 5 lists all land use codes in the dataset in Section 4.1. The original description in French is given, along with a translation.

Code	Original description	English translation
AB	terrains à bâtir	building plot
AG	terrains d'agrément	amenity land
B	bois	woodland
BF	futaies feuillues	deciduous forest
BM	futaies mixtes	mixed forest
BO	oseraies	willow grove
BP	peupleraies	poplar grove
BR	futaies résineuses	coniferous forest
BS	taillis sous futaie	coppice with standards
BT	taillis simples	simple coppice
CA	carrières	quarry
CH	chemin de fer	railway
E	eaux	water body
J	jardins	garden plot
L	landes	heathland
LB	landes boisées	wooded heath
P	prés	meadow
PA	pâtures	pasture
PC	pacages	grazing land
PE	prés d'embouche	grazing meadow
PH	herbages	grassland
PP	prés plantes	crop meadow
S	sols	residential plot
T	terres	arable land
TP	terres plantées	planted arable land
VE	vergers	orchard
VI	vignes	vineyard

Table 5: Land-use classification codes for the variable `code_nature_culture` in the dataset introduced in Section 4.1. The original French description is provided alongside the English translation.