

Review

# MACHINE LEARNING OF SPATIAL DATA

Behnam Nikparvar<sup>1,†,\*</sup>  and Jean-Claude Thill<sup>2,3,†,\*</sup> 

<sup>1</sup> Infrastructure and Environmental Systems Program, the University of North Carolina at Charlotte, 9201 University City Blvd., Charlotte, NC 28223, USA; bnikparv@uncc.edu

<sup>2</sup> Department of Geography and Earth Sciences, University of North Carolina at Charlotte, 9201 University City Blvd, NC 28223, USA; jfthill@uncc.edu

<sup>3</sup> School of Data Science, University of North Carolina at Charlotte, 9201 University City Blvd, NC 28223, USA; jfthill@uncc.edu

\* Correspondence: bnikparvar@gmail.com; jfthill@uncc.edu

† These authors contributed equally to this work.

**Abstract:** Properties of spatially explicit data are often ignored or inadequately handled in machine learning for spatial domains of application. At the same time, resources that would identify these properties and investigate their influence and methods to handle them in machine learning applications are lagging behind. In this survey of the literature, we seek to identify and discuss spatial properties of data that influence the performance of machine learning. We review some of the best practices in handling such properties in spatial domains and discuss their advantages and disadvantages. We recognize two broad strands in this literature. In the first, the properties of spatial data are developed in the spatial observation matrix without amending the substance of the learning algorithm; in the other, spatial data properties are handled in the learning algorithm itself. While the latter have been far less explored, we argue they offer the most promising prospects for the future of spatial machine learning.

**Keywords:** spatial machine learning; spatial dependence; spatial heterogeneity; scale; spatial observation matrix; learning algorithm; deep learning

## 1. Introduction

Machine Learning (ML) has become a widely used approach in almost every discipline to solve a broad range of tasks and problems with structured and unstructured data, including but not limited to regression, grouping, classification, and prediction. It has proved itself to be a powerful and effective tool in various disciplinary fields and domains of application where spatial aspects are essential, including the following: land use and land cover classification [1,2], cross-sectional characterization [3,4] and longitudinal change [5], urban growth [6] and gentrification [7], disaster management [8], agriculture and crop yield prediction [9], infectious disease emergence and spread [10], transportation and crash analysis [11], map visualization and cartography [12,13], delineation of geographic regions [14] and habitat mapping [15], geographic information retrieval and text matching [16], POI and region recommendation [17], trajectory and movement pattern prediction [18], point cloud classification [19], spatial interaction [20], spatial interpolation [21], and spatiotemporal prediction [22–24].

Spatial data exhibit certain distinctive properties that set them apart from other data types, such as spatial dependence, spatial heterogeneity, and scale. As in other modeling approaches, we need to be aware of the specificities that these properties entail when we conduct ML on spatial data. Indeed, the explicit handling of these spatial properties can improve the performance of the ML model or add meaningful insights into the process of learning a task. At the same time, failure to appropriately include these properties into the ML model can negatively impact learning. At this juncture, there is extensive literature that applies ML to spatial data but research that explicitly features the spatial

properties of data in ML remains rather limited. Therefore, a survey of this body of research and of their findings is crucially needed to fully understand how much progress in handling spatial data in ML has been accomplished, what the best practices in this respect are, what gaps may remain in this literature, and where opportunities exist for future research with, and on, spatially explicit ML.

Surveys of the extant literature have previously been conducted on several themes intersecting with spatial ML, namely knowledge discovery and data mining [25–27], spatial prediction methods [28], artificial neural networks in geospatial analysis [29], ML for spatial environmental data [30], and active deep learning for geo-text and image classification [31]. While their contribution to furthering scholarship in spatial ML is unquestionable, these surveys are rather limited in scope. For instance, while there are overlaps between data mining and ML, they have distinct definitions, follow different processes, and have different goals. The emphasis of other literature surveys is on a specific discipline or learning method. In addition, the focus has usually been on applications, and a detailed discussion of spatial properties and their role in the ML process are still missing. Thus, this paper aims to review some of the best recent practices of ML of spatial data while considering the learning process and the properties of spatial data. Interested readers can explore more deeply the topics broached in this paper by following some of the references provided herein. Furthermore, our intention is not to provide a general comparison between ML methods. For such purpose, interested readers are referred to [5].

The rest of this paper is organized as follows. We begin with a brief overview of ML and of spatial properties of data in the following two sections. Then we will lay out the process of ML of spatial data in two steps entailing the construction of a spatial observation matrix and a learning algorithm. The following section is concerned with the ML of spatiotemporal data. Finally, we summarize the review work conducted for this paper, discuss the gaps in the current state of knowledge and practice in ML for spatial data and identify possible areas of fruitful research for the future.

## 2. Machine Learning

ML can broadly be defined as the capability of a computer program to improve automatically with experience via the performance of certain tasks. A performance measure quantifies this experience, and if it improves, we say the machine is learning [32]. As shown in Figure 1, ML can be classified into three types: unsupervised, supervised, and reinforcement learning.

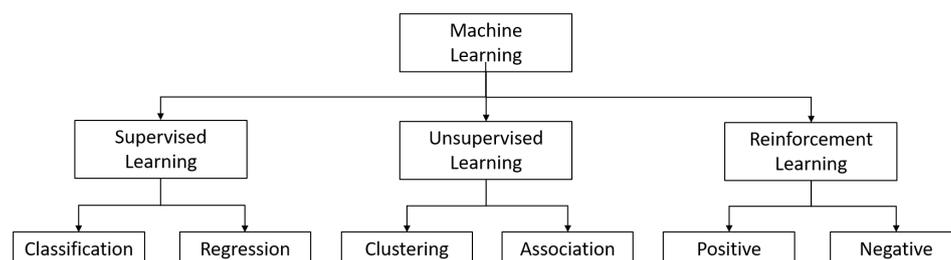


Figure 1. Machine learning.

In learning tasks, we usually aim to estimate one or more output variables  $Y = (Y_1, \dots, Y_m)$  for a given set of input variables  $X = (X_1, \dots, X_n)$ . When the desired output variables  $Y$  are in hand, the learning task is dubbed supervised learning or metaphorically learning from a teacher [33]. In other words, we know the correct answer, and we try to learn the dependency between input variables  $X$  and  $Y$ . At the same time, one can take  $(X, Y)$  as random variables since many factors may influence data and measurements, making the whole setting stochastic [30]. Then, we can represent these random variables by a joint probability density  $P(X, Y)$ . In this case, a supervised learning is concerned with determining the properties of the conditional density function  $P(Y|X)$ .

Output variables  $Y$  could be class labels (classification) or real numbers (regression), possibly resulting from the coding of unstructured data [34].

When the desired target variables  $Y$  are not obtained, the learning task is called unsupervised learning or learning without a teacher [33]. In this case, the training only involves the  $n$  observations of random variables  $X$  with joint density  $P(X)$  and the goal is to directly infer the properties of this probability density function. These properties can be finding the joint values of the variables  $X = (X_1, \dots, X_n)$  that frequently appear in the data (association rules), grouping or segmenting a collection of similar or dissimilar samples into subsets (clustering), or projection of data into lower dimensions usually ranked by variance (dimensionality reduction). Unsupervised learning is typically concerned with finding patterns in the data and thus close to data mining and knowledge discovery in databases [34].

Existing somewhere between supervised and unsupervised learning, semi-supervised learning is based on extending either unsupervised or supervised learning to include additional information typical of the other paradigm. Two main settings of semi-supervised learning are semi-supervised classification and constrained clustering. The former is a classification task with partially labeled data (usually useful when the training sample size is small). The latter is unsupervised learning with some sort of supervised information about the clusters [35]. Recently popular, active learning methods are a subset of semi-supervised machine learning. The idea behind active learning is that a machine can learn with less training data if it is allowed to choose from the training data set by asking questions [36]. A comprehensive review of these methods for application in (geo) text and image classification is available in [31].

Finally, in reinforcement learning, the machine produces some actions and interacts with its environment. These actions affect the state of the environment and, in turn, receives some scalar rewards (positive reinforcement learning) or punishments (negative reinforcement learning). The goal of the machine is to learn to act in a way that maximizes (minimizes) the future rewards (punishments)[34].

### 3. Spatial Data Properties

In machine learning, observations are represented by a matrix  $X$  where the rows are instances (samples) of a phenomenon under study, and columns are different attributes associated with each of these instances. The same applies to spatial data, but the samples on each row are also referenced to a specific location in the geographic space. To define the relationship between the real world and this matrix, we can choose between two well-known views of the world, namely the field view or the object view [37]. Field entities are usually represented by regular grids and object entities by points, lines, or polygons. Being referenced to a specific location in space creates unique properties for spatial data that geographers and econometricians have studied over the years [38–40]. There is a broad agreement in the literature that there are three fundamental properties for spatial data: spatial dependence, spatial heterogeneity, and scale.

#### 3.1. Spatial Dependence

Named as the first law of geography by [41], "near things are more related than distant things" formulates the first property of spatial data, which is known as spatial dependence. Spatial dependence is a fundamental and useful property of spatial data that stems from the general continuity of space. A variety of ways exist to express and measure spatial dependence in data sets. For example, a spatial autocorrelation statistic can summarize the similarity of the values for a variable of interest at different locations as a function of the distance that separates them or of their adjacency to each other [42,43]. The Global Moran's  $I$  is an indicator of spatial autocorrelation and is a similarity index, which is usually used for areal data and is calculated based on the cross-product of variations from the mean [44]:

$$I = \frac{\sum_i \sum_j W_{S_i S_j} (Z(S_i) - \bar{Z})(Z(S_j) - \bar{Z})}{\sigma^2 (\sum_i \sum_j W_{ij})} \quad (1)$$

Where  $Z(S_i)$  and  $Z(S_j)$  are the attribute values at sample locations  $S_i$  and  $S_j$ , respectively.  $\bar{Z}$  is the sample mean and  $\sigma^2$  is the sample variance. Weight matrix element  $W_{S_i S_j}$  represents the spatial relationship between paired geographic units, which can be defined based on contiguity (topological relationship of spatial units) or distance (distance is determined on a physical or social network, 2D or 3D Cartesian space). The  $I$  index takes positive or negative values between -1 and +1. A positive value means similar values happen in close proximity, while a negative value shows dissimilar values are spatially grouped close to each other. A value close to zero shows no spatial autocorrelation for the variable, indicative of a spatially random process; in the latter case, the assumption of independence essential for many statistical methods is met. More measures of spatial autocorrelation have also been suggested, such as Geary's  $C$  [45], which considers the square of variations of the variable between two locations. Alternatively, semivariogram analysis (see [46]) can be conducted to estimate the spatial autocorrelation structure of a stochastic process based on the variance of the observations over a range of distances (usually for geographic phenomena represented by points). The empirical semi-variogram  $\hat{\gamma}(h)$  is defined as half the average squared difference between values of a random field  $Z$  at pair points  $S_i$  and  $S_j$  for a given distance  $h$  in a region:

$$\hat{\gamma}(h) = \frac{1}{2d(h)} \sum_{|S_i - S_j|=h} (Z(S_i) - Z(S_j))^2 \quad (2)$$

The semi-variogram is visually analyzed in a variogram. A range parameter is defined as the distance beyond which spatial autocorrelation disappears in the data, thus expressing the spatial dependence structure. In principle, spatial autocorrelation decreases as the distance among geographic units increases [47]. From a statistical point of view, spatial dependence violates the assumption of independence and identical distribution (i.i.d.), central to many statistical methods. While this assumption is not required for ML methods, spatial dependence properties still need to be considered for some applications and may significantly contribute to enhance the quality of the learning process, which will be discussed later (see section 4.1.1).

Two more complex components of spatial dependence are the neighborhood effect and spillover effect. While the spatial weight matrix is usually used to capture the spatial dependence of the dependent variable, spatial association rules, which are based on linguistic or topological (in mathematics language) rules, are used to capture these two more implicit spatial effects. Association rules such as: "if block group A is next to a high crime neighborhood, then block group A has high crime" (neighborhood effect) or "if a block group A is next to a shopping mall, then block group A will experience high crime" (spillover effect) are examples of these implicit spatial effects [48]. However, most traditional machine learning algorithms do not consider the impact that spatial dependence property may have on learning these association rules. For example, while we expect the crime rate at a block to be more similar to its neighboring blocks than the farther blocks, traditional machine learning algorithms like Support Vector Machines (SVM) or Neural Networks have not been designed to incorporate this characteristic [30].

### 3.2. Spatial Heterogeneity

All patterns that we can observe result from the four main processes, namely interaction, dispersion, diffusion, and exchange over the geographic space [37]. A specific observed pattern can be the result of one or multiple processes. These processes may happen in different spatial as well as temporal scales, over varied durations, and

with differentiated intensities. A mixture of those with the place-based environmental factors and contexts (which are the existing patterns resulted from interweaving previous processes at that location) in different sub-regions creates even more complex patterns of outcomes [49,50].

Global measures of spatial autocorrelation may confirm the existence of positive or negative self-similarity with regard to distance, but this comes at the cost of a fundamental assumption. The parameters (mean and variance) of the random function representing the process are assumed to be constant, which means that the sample's distribution is even over the extent of the territory over which data are generated. This is called stationarity of the random function associated with that process [46], and when it is violated (called a nonstationary process), the process is heterogeneous. In other words, a spatial process is said to be stationary when the difference between values of an attribute is only explained by the distance between the points or units [51,52]. Another source of spatial heterogeneity is when the spatial dependence is different in various directions (anisotropy) For example, high precipitation patterns may be interrupted in a specific direction where the spatial topography of the terrain (mountains) blocks the clouds [28].

### 3.3. Scale

Scale is also important because it can inform about sampling for training experience. Learning is more reliable when the distribution of the samples in the training experience is similar to the distribution of the test experience [32]. In many geographic studies, training occurs on data from a specific geographic area. This makes it challenging to use the trained model for other geographic regions because the distribution of the test and train datasets is not similar due to spatial heterogeneity [53]. This means the sampling strategy for the training dataset is essential to cover the heterogeneity of the phenomena of interest over the spatial frame of study. Scale can inform the sampling of training dataset by defining its elements; namely, resolution (measurement scale), context (scale at which the process is operating), and spatial extent (extent of observation) [54]. By increasing the extent of the study area, more processes and contextual environmental factors may alter the variable and result in non-stationarity by interweaving spatial patterns of different scales or inconsistent effect of processes in different regions [55]. This is especially important because collecting spatial samples is an expensive undertaking.

There are two distinct challenges related to scale and zoning when working with spatial data with areal units: the modifiable areal unit problem (MAUP) and the uncertain geographic context problem (UGCoP). The first notion is related to the sensitivity of the analytical results to the definition of the geographic units for which data are collected [56]. [57] demonstrated how different levels of spatial aggregation and zoning could result in different values of the correlation coefficient for areal units. They generated two variables with high positive and negative spatial autocorrelation and investigated the effect of varying levels of aggregation in the correlation coefficient. Results show that even in the absence of spatial autocorrelation, the correlation coefficient increases by grouping and aggregation. Attempts have been made to provide solutions for the MAUP problem based on the size and interconnectedness of the areas [58] and spatial entropy [59,60].

The second problem, UGCoP, refers to the sensitivity of the contextual variables and analytical results to different delineations of contextual units [61]. Kwan highlights how varying delineations of contextual units, even if everything else is the same, can lead to uncertain and inconsistent analytical results over time. Notice this problem is different from the MAUP because it refers to the contextual influences on the individuals being studied in geographic units with unknown or uncertain spatial configuration. In this regard, it is closer to the ecological fallacy [62] problem, which relates to inferences about individuals from inferences in the aggregated level. An example is in environmental health studies, where residential neighborhoods (e.g., census, postal code, or buffer) are

typically used as contextual units. However, these geographic units may not accurately represent the actual areas that apply contextual influences on the health outcome. People are exposed to health risks in different locations (home, school, etc.) during the day, and it is not easy to delineate the boundaries of such exposures. UGCoP is not due to different zoning or spatial scale, but contextual influences naturally change across granularities, making it even more complicated to delineate geographic contextual units [61].

The scale of analysis and appropriate geographic contextual units are among the first questions that one may need to answer before applying any spatial machine learning. However, the above-mentioned fundamental problems have rarely been investigated in spatial machine learning applications. For example, [63] shows how MAUP can lead to perturbations in the convolution-based residual neural networks used for urban traffic prediction. Thus, the investigation of such effects on spatial machine learning is critically needed.

### 3.4. Other Properties

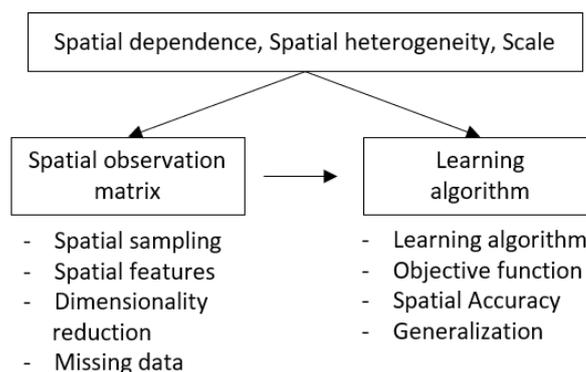
In addition to the above main properties of spatial data, some less fundamental but still essential properties may result from the specific representation of the data (polygon, line, point, regular grid, text) or the process of measurement [37]. For example, discretizing the continuous space into pixels may cause a loss of information at the sub-pixel level [64], while delineating the boundaries of geographic regions entails some generalization and may not be easy due to the Boolean nature of typical regionalization approaches [65]. When working with text data, disambiguation in place names and addresses due to multiple instances for a single entity, a single entity with multiple names, and different addressing systems can create uncertainties [14,16,66,67]. Also, the geometric interpretation of vague place names like 'midland' or 'near' is usually not straightforward [68,69].

Moving from a 2D spatial space to a 3D space may impose some limitations and create some biases when ML applications are used. An example would be the case of 3D point clouds. The density of the point cloud usually changes by distance from the sensors. The closer the sensor, the higher the density. The implication of this effect is that identifiable spatial features in more dense areas of the cloud may not be recognizable in more sparse areas [70].

Finally, for many image-based applications (less for satellite images and more for close-range images), the orientation of the sensor at the time of capturing the image may influence the amount of energy received by the sensor. This makes orientation an important element. A case susceptible to such effect in remote sensing would be that of night light images, where satellite viewing vertical and zenith angles can significantly impact the amount of light in urban areas [71]. Statistical models are commonly used to handle these characteristics in remote sensing. In close-range photogrammetry and computer vision applications, however, the sensor's geometry at the time the image is captured is usually reconstructed. Interested readers are referred to the famous structure from motion [72] and visual simultaneous localization and mapping (Visual SLAM) [73]. As for applications in ML, the simplest way to make the model invariant of the orientation is to train the model with samples from different orientations and image augmentation.

## 4. Machine Learning of Spatial Data

To conduct machine learning of spatial data, we need to add location, distance, or topological relations to the process of learning. Figure 2 organizes the learning process into two steps, the spatial observation matrix and the learning algorithm. We reviewed the literature to address how spatial data properties can be involved in each of these steps.



**Figure 2.** Machine learning of spatial Data.

#### 4.1. Spatial Observation Matrix

One typical way to include spatial properties in ML is to find a representation for these properties in the observation matrix  $X$ . The principle here is that, after we design and engineer the observation matrix  $X$  to include spatial properties, we can effectively use typical ML methods (e.g., families of decision trees and random forests, support vector machines, neural networks, ensemble models, etc.) without making any change to the learning algorithm. Several critical aspects are involved in creating an ultimate spatial observation matrix used as an input to the learning algorithm, namely spatial sampling, spatial features, dimensionality reduction, and handling missing data. These are discussed hereunder.

##### 4.1.1. Spatial Sampling

While tremendous progress has been made in spatial data collection technologies, ML methods still face essential challenges in acquiring optimized samples for training. The current view in ML is to move from model-centric ML to a more data-centric ML [74]. From a statistical point of view, the size and distribution of a sample set should represent the entire population or distribution. Two important points need to be made here. First, this is distinct from the commonplace concern about the poor arrangement of samples in training and test data sets in ML, which leads to a typical generalizability issue (See section 4.2.11 on this matter). Instead, the entire sample set (including both training and test data sets) should represent the phenomenon being learned, which is a sampling problem [75].

Second, the representativeness of the samples is defined in the attribute, spatial, or temporal space (or multiple spaces at a time) depending on the application. The structure of the data in each space may impose some special properties. As far as the spatial data matrix is concerned, data properties in a spatial space, such as spatial dependence, can inform an optimized sampling configuration to avoid redundancy.

It is not always the scarcity of samples that leads to challenges for learning. Over-sampling will not impact the learning process because the assumption of i.i.d. is not required in ML [76]. However, it may overestimate the accuracy of learning in the assessment process. For example, let us consider a land cover classification task using satellite imagery. Suppose a large batch of samples are selected from close locations for a single class category (e.g., vegetation). It is most likely that these samples are similar to each other. In this case, even if samples from other sites are available in the sample data set, the classifier is somehow biased since it most probably labels the frequent familiar samples correctly. Thus, the confusion matrix class accuracy and recall, which are usually used to evaluate classification results, are not robust to assess the learning results with the varying sample sets [77,78]. This problem is known as intra-class imbalance [79]. Intra-class imbalance is not solved by cross-validation methods typically used in the generalization step (Section 4.2.11). That means we either need an evaluation method that can account for such effects or be careful in sampling the data that enter the learning

algorithm. For a useful review of the spatial sampling methods, see [47]. It is worth noting that inter-class imbalance, in which the number of samples in class categories are highly uneven, can also degrade the accuracy of classification. The performance on a class with many samples will be higher, compared to the classes with fewer samples. While it overestimates the overall accuracy, inter-class accuracy is still identifiable by looking at single class accuracy and recall in the confusion matrix [80].

#### 4.1.2. Spatial Features

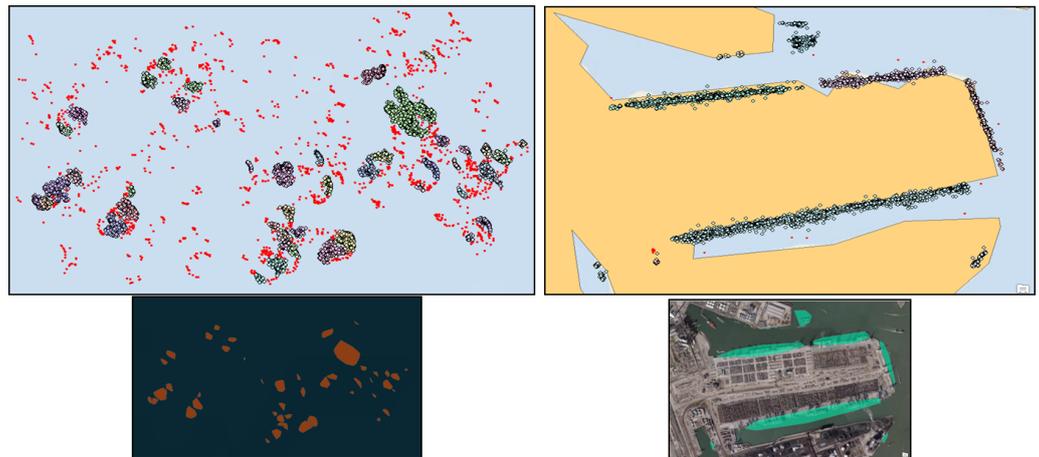
Several methods exist to include the spatial components of data into the observation matrix. One way is to directly add the spatial reference to the data matrix as attributes. This entails embedding the spatial references of data directly into the attribute space. Practically, this can be implemented in either of two ways. One consists in adding coordinates (e.g., latitude and longitude) alongside semantic attributes for each observation to the observation matrix [81,82]. However, using the coordinates as spatial variables may generate considerable overfitting because they are highly correlated [83]. The other typical way is to add observations tied to a region as fixed effects of that region to the observation matrix [84,85]. This approach is effective to handle inclusion relationships, it cannot capture complex structures, however.

In addition to spatial reference information, spatial entities and phenomena have within object information (geometric, spectral, textural, and statistical) and between-object information (contextual and relational) that can be created as new features and added to the observation matrix directly [1]. Geometric information of single spatial entities can be used in different forms such as length, area, and ratio thereof. Spectral and textural information have extensively been used in the remote sensing community for land-cover and land-use classification. However, spectral features are insufficient for this purpose due to heterogeneity, especially in urban environments [2].

Other features such as texture, which indicate coarseness or smoothness of a pattern in an image, have been suggested as complementary information. From a mathematical point of view, numerous functions exist that can represent the local variability of pixel values as texture features, including first-order statistics (e.g., means, variance, and standard deviation) and second-order statistics (gray-level co-occurrence matrix, spatial autocorrelation) [86]. These methods, however, focus on a single scale. Multi-resolution spatial and frequency analysis tools, such as Gabor transform, Wigner distribution, and wavelet transforms, have effectively been used to overcome this problem [87–90].

Texture analysis is usually conducted using moving windows (regularly shaped grid) or fields (irregular shape). Compared to moving windows, fields partition the area into homogeneous regions and provide a more realistic representation of the spatial entities [86]. The boundaries of these fields are determined using existing polygon features or digital image segmentation methods (e.g., edge detection, region growing) [91]. In addition to textural and spectral information, geometric and contextual attributes of the fields can also be used. In this respect, [92] used the proportion of build-up area, vegetation, and water surfaces, and [93] calculated spatial metrics for land covers and semivariogram for Normalized Difference Vegetation Index (NDVI).

We provide here another example of how to identify and delineate the boundaries of functional zones based on moving object patterns. Figure 3 shows two zones inside and outside the port of Rotterdam, the Netherlands. The points represent the recorded location of vessels that visited an unlabeled zone and are extracted from vessel trajectories stop and move segments. On the one hand, the point clusters may inform about the boundaries of each zone. On the other hand, shape and simple statistics within each zone such as the duration of visits, the number of unique vessels, distribution of vessel course (direction), and vessel type may provide some information about the functionality of space (e.g., anchorage, containerized ship berth, bulk ship berth, etc.). These within entity information can be added as new features to the observation matrix to recognize the functionality of the zones.



**Figure 3.** Figure 3. Anchorage regions (left) and containerized berth regions (right) outside and inside the port of Rotterdam. Points represent the vessel trajectory instances every one hour. Point clusters may inform about the boundaries of each zone (bottom). In addition, shape, the duration of visits, number of unique vessels, and vessel types within each region can inform its functionality.

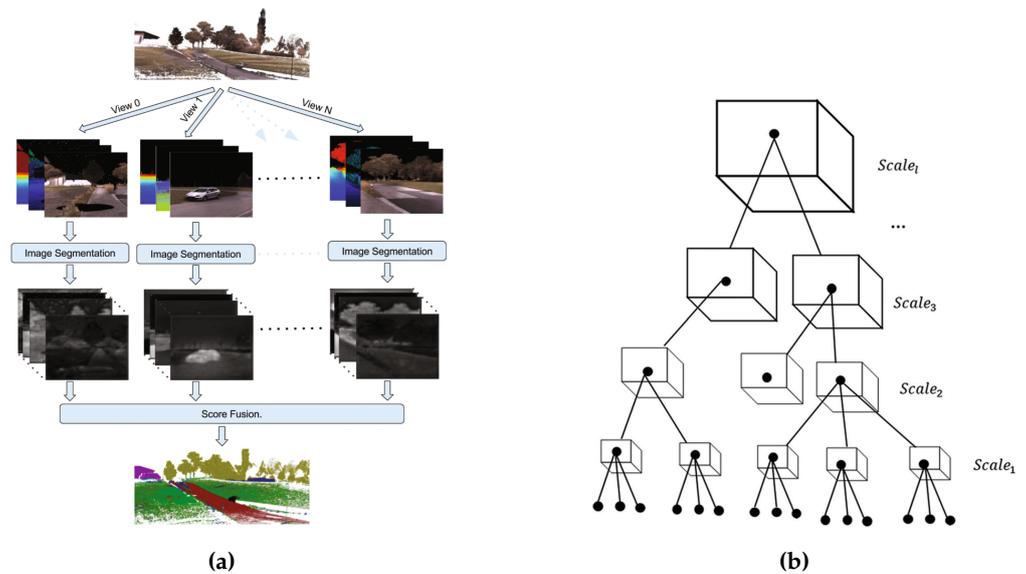
Further, between-objects information such as connectivity, contiguity, distance, association, and direction can be used to create new features [1]. In the vessel movement example, the association of zones where different types of ships dock to load and unload cargoes in close distance may become a discriminating feature to distinguish the port area from anchorage zones outside of the port (in Figure 3, compare bottom-left and bottom-right quadrants).

Point cloud classification is another area where machine learning methods have been used extensively. The simplest way to classify point clouds is to add the third coordinate as an attribute of each point and use standard 2D ML methods to classify the point cloud. However, there are two critical shortcomings for such approaches [70]. First, point clouds may have multiple  $z$  values for a point with  $(x,y)$  coordinates, and compression of the 3D point cloud to 2D image causes a loss of information. Second, points in point clouds are usually irregularly spaced, making the selection of a fixed window size difficult. In other words, the point densities may vary in relation to the distance to the sensor. As a result, features learned in a dense area are not generalizable to sparse areas.

[94] used regularly spaced voxels and voxel feature encoding (VFE) to address the first problem. This method is subject to loss of information due to decreased spatial resolution and increased memory usage in the voxelization process. Alternatively, [95] projected point clouds on multiple synthetic 2D images and labeled pixels based on the prediction scores from these synthetic images to handle the first issue (Figure 4a). [70] suggested creating a point cloud pyramid with  $I$  scale levels by subsampling the original point cloud. The deep hierarchical features (see section 4.2.6) for each point are extracted using a deep neural network within each scale. Such an approach forms a feature pyramid (Figure 4b). The feature vectors of a point along this pyramid are concatenated to create a final feature vector that is fed into a classifier. This final feature vector contains both hierarchical and multi-scale features of the original point cloud, which can address both issues discussed earlier.

#### 4.1.3. Dimensionality Reduction

ML tasks may end up with a large number of input variables in the observation data matrix. A disproportionately large number of interrelated variables may negatively impact learning in a variety of ways. Apart from the need for more training data and an increase in processing time, an unnecessarily large number of variables may impose new uncertainties into the learning process because of the correlation between variables [96].



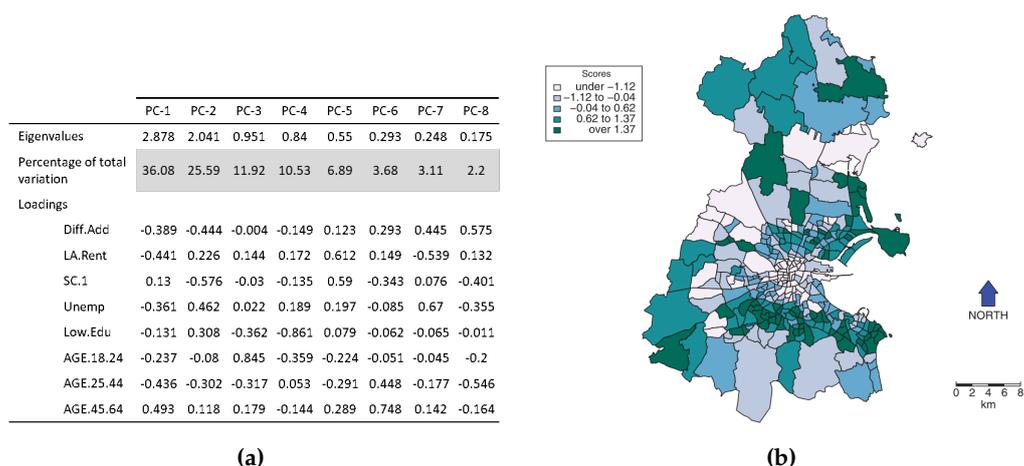
**Figure 4.** (a) Method proposed by [95] (Reprinted by permission from Springer Nature: Springer, [Deep projective 3D semantic segmentation](#) by Lawin, F.J.; Danelljan, M.; Tosteberg, P.; Bhat, G.; Khan, F.S.; Felsberg, M. COPYRIGHT (2017)). The input point cloud is projected into multiple virtual camera views, generating 2D images. The images for each view are processed for semantic segmentation. The output prediction scores from all views are fused into a single prediction for each point, resulting in a 3D semantic segmentation of the point cloud. (b) The multi-scale representation of point clouds in [70] ([Using multi-scale and hierarchical deep convolutional features for 3D semantic classification of TLS point clouds](#). Guo, Z.; Feng, C.C. International Journal of Geographical Information Science 2020,34, 661–680., Taylor Francis. reprinted by permission of the publisher), with the finest  $scale_1$  at the bottom and the coarser  $scale_1$  on the top.

One way to handle such a problem is to select a subset of features that provide the best results. A variety of methods exist for feature selection, such as a genetic algorithm [97]. Dimensionality reduction methods are another way to handle many useful variables, especially when the influence of each variable is not of interest.

To reduce unnecessary variables, we need to understand the structure of the variance-covariance matrix. The problem is calculating the variance-covariance matrix  $C_{n \times n}$  for a given observation matrix  $X_{m \times n}$  is computationally expensive for a large number of variables. Several dimensionality reduction methods exist, including but not limited to Principal components analysis (PCA), factor analysis, independent components analysis, and self-organizing maps (SOM) [98–100]. PCA is a statistical method well known for its application in dimensionality reduction. PCA transforms the observation matrix  $X_{m \times n}$  with many interrelated variables to a smaller set of new uncorrelated variables (components). These new variables are ordered so that the first few retain most of the variation present in the original variables [99].

It is proved that by applying useful linear algebra operations,  $C$  can be decomposed into  $LVL^T = C$ , where  $L$  is a matrix of eigenvectors (representing the weight of each variable on corresponding principal components), and  $V$  is an orthogonal matrix of eigenvalues (representing the variances of the corresponding principal components). The observations are then transferred to a space where the column vectors in  $L$  represent the axes. In practice, a few principal components that represent most of the variation in the data set are selected (Figure 5). These transformed observations are then entered into the learning algorithm.

For an observation matrix with spatial references, a simple PCA assumes the structure of the variance-covariance matrix remains stationary across the study area. [101] showed that PCA could be weighted locally either in attribute space (LWPCA) or in geographic space (GWPCA) to account for certain heterogeneities. For the former, we



**Figure 5.** [101] (Geographically weighted principal components analysis. Harris, P.; Brunson, C.; Charlton, M. International Journal of Geographical Information Science 2011, 25, 1717–1736, Taylor Francis. reprinted by permission of the publisher ). (a) Summary of global PCA for eight variables representing social structure in Greater Dublin (the first three components express more than 70% of variations). (b) Spatial distribution of first component scores from global PCA.

assume the covariance structure is homogeneous for observations that are close to one another in attribute space, which leads to  $L_i V_i L_i^T = C_i$  with respect to the sub-region  $i$ . The GWPCA for location  $(S_i, S_j)$ , however, is written as  $LVL^T|(S_i, S_j) = C(S_i, S_j)$ , with  $C(S_i, S_j)$  as the variance-covariance matrix of location  $(S_i, S_j)$ .

While PCA is a common approach in summarizing (sub)sets of variables in ML, LW-PCA and GWPCA have not been investigated. Such methods can additionally provide insights into the spatial distribution of each composite variable by mapping components scores. This may lead to a complete understanding of a process [102]. Moreover, eigenvalues and eigenvectors can be obtained for locations with no observations. Challenges include bandwidth selection and more computational cost. A locally weighted PCA can also optimize sampling re-design by identifying local and spatial outliers rather than global and aspatial outliers[101].

#### 4.1.4. Missing Data

Nowadays, data are more available, both spatially and temporally. However, given they are more often organic, being a byproduct of the processes that created them in the first place, there are always gaps in temporal and spatial dimensions because of effects and circumstances that are out of our control. Therefore, missing data is an important challenge, and many analyses are simply not implementable without having a way to deal with this problem. Missing observations can be independent of each other, dependent on their neighboring points, or with specific patterns [103]. There are different ways to address missing values, such as aggregating data into coarser granularity, removing the observations with missing values from the data set, and imputing values. Although imputing data adds preprocessing step to analysis, it leverages the existing data and avoids losing information due to aggregation and discarding some observations.

Spatial prediction methods can always be used to impute values for data sets with missing values. The most famous approaches for spatial prediction are spatial statistical models (e.g., geographically weighted regression) and geostatistical models such as kriging [104]. Several studies have demonstrated that kriging in the universal form is preferred to geographically weighted regression for prediction due to its optimal statistical properties [105].

Apart from the statistical and geostatistical models, other approaches in machine learning and other paradigms have also been adapted to spatial properties for missing

data imputation. Probabilistic Principal Component Analysis (PPCA), for example, is a probabilistic extension of PCA and has been used to impute missing values for different applications [103,106]. This method has proven useful when a significant portion of the observation matrix is unknown [107]. A comparison of different types of PPCA can be found in [108], and an extension of PPCA to consider temporal and spatial dependence can be found in [109]. All systematic errors must be removed before using PPCA for imputation [103].

#### 4.2. Learning Algorithm

Instead of generating new spatial features and process them with traditional non-spatial machine learning methods, we can directly incorporate spatial properties in the learning algorithm. Among all ML methods, decision trees and random forests, support vector machines (SVM), neural networks, and deep neural networks (DNN) have found considerable attention in spatial science.

##### 4.2.1. Decision Trees

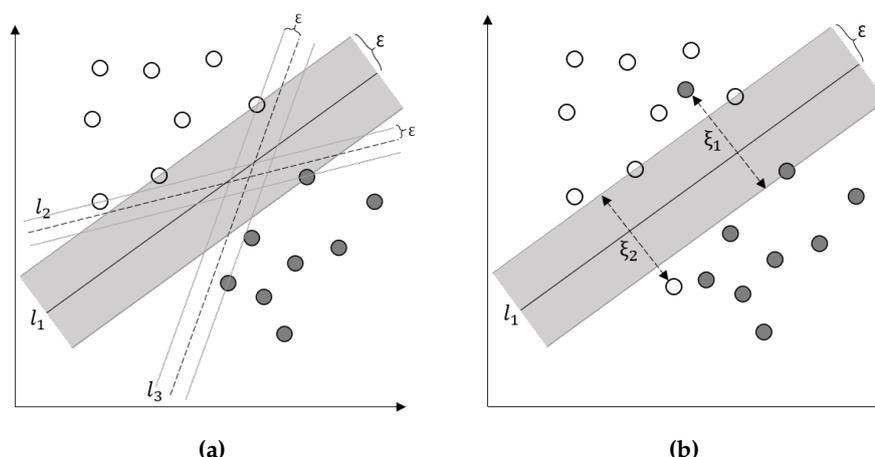
Decision trees (DT) are popular ML methods adapted for spatial problems to overcome violation of the i.i.d assumption. As a class, spatial entropy-based decision tree classifiers use information gain coupled with spatial autocorrelation to select candidate tree node tests in a raster spatial framework. For example, [110] added a spatial autocorrelation measure to the target function evaluated at each node of the tree. PCT is a multi-task approach where hierarchies of clusters of similar data are identified, and a predictive model is associated with each group. When splitting a group is considered at a node, a test is run that maximizes within-cluster variance reduction. To account for spatial non-stationarity in the target variable, a term based on global measures of spatial autocorrelation was added to this test.

A common trait of such approaches is to add a constraint to account for spatial properties, while still relying on local entropy testing at tree nodes. One of the frequently occurring issues in image classification using decision trees is salt and pepper noise. Salt and pepper noise happens when the predicted label of a specific pixel differs from its neighborhood pixels and can result from high spatial autocorrelation in class labels of the sample data used for training. [111] proposed a focal-test-based spatial decision tree (FTSDT), in which the tree traversal direction of a learning sample is based on local and focal (neighborhood) properties of features. They use local indicators of spatial association-Lisa [112] as spatial autocorrelation statistics to measure the spatial dependence between neighborhood pixels.

##### 4.2.2. Support Vector Machines

Support vectors machines (SVMs) have been used for classification and regression problems [113]. The idea of SVM is to map the original input space to a higher-dimensionality feature space where the observations are separable by hyperplanes (Figure 7a). Among all possible hyperplanes, the one that maximizes the margin width ( $\epsilon$ ) is optimized ( $l_1$ ). In reality, observations may not be separable easily due to the presence of outliers (Figure 7b). Instead of increasing the complexity of the model structure (in this example, a nonlinear curve), we allow misclassification of some observations and penalize them based on their distance from the margins ( $\xi$ ). In essence, the objective function has two terms: a term containing  $\epsilon$  and another term containing  $\xi$ , and the goal is to maximize the former and minimize the latter.  $\xi$  is called a regularization term. Regularization terms are usually added to objective functions to control the complexity of the model and avoid overfitting (see section 4.2.11). SVM performs well in high dimensional spaces. It is less sensitive to class imbalance and powerful in generalization [11].

[114] suggested an extension of SVM called support vector random field that explicitly models spatial dependencies in the classification using conditional random fields



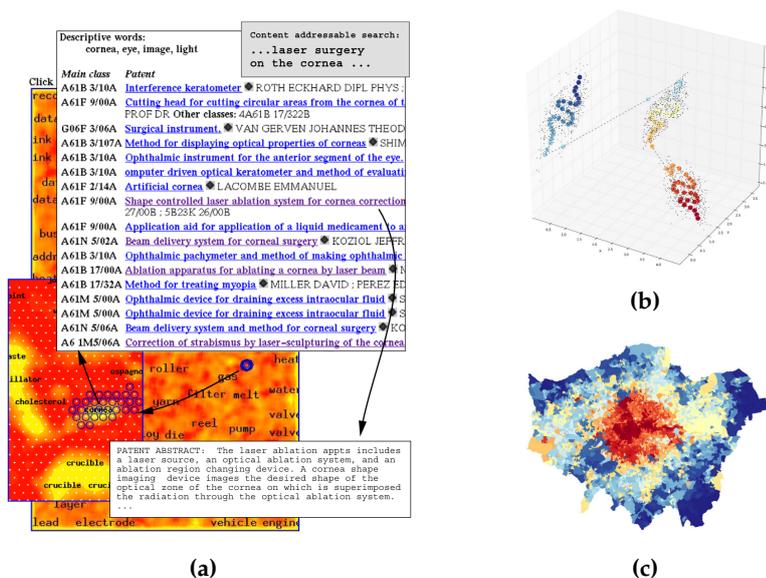
**Figure 6.** Observations are transformed to a higher-dimensionality space (for simplicity, a 2D space). (a) Among all possible hyperplanes (line or curves in the current 2D example), the hyperplane that maximizes the margin  $\epsilon$  is optimum ( $l_1$ ). (b) In practice, outliers exist in observations. While there is still a hyperplane that can separate the two classes, it may result in overfitting.

(CRF). The model contains two components: the observation-matching potential function and the local-consistency potential function. The former models the relationship between the observations and the class labels using an SVM classifier, and the latter models the relationship to neighborhood labels. The local-consistency potential function penalizes discontinuity in between the pairwise sites.

#### 4.2.3. Self-Organizing Maps

Self-Organizing Maps (SOM) are among nonlinear clustering methods that have been used with spatial and non-spatial data [100]. SOM is a simple neural network with no hidden layers. It maps an  $n$  dimensional feature vector to a regular grid of square (four neighbors) or hexagons (six neighbors) neurons in the output layer, initialized with  $n$  weights. First, a similarity measure is used to find more similar neurons to the input feature vector. Then, the weights of the activated neurons and their neighboring neurons are adjusted to make them even more similar to the input vector. This process is repeated for a set of input feature vectors. Finally, it creates a spatial organization of the neurons in a one-, two-, or three-dimensional space in which dissimilar units stay farther away. SOM is similar to K-means clustering but different in two ways. First, K-means is based on the nearest distance, while SOM utilizes distances between all paired neurons (weighted by a neighborhood kernel) [20]. Second, SOM also visualizes the relation between clusters by representing how far they are from each other in a topological space [100]. Such property makes SOM attractive for visual data mining. For example, it is possible to compare the SOM visualizations with other forms of visualization (e.g., geographic visualization), especially in an interactive platform [115]. [12] used SOM to visualize demographic trajectories, and [20] demonstrated how SOM could be employed to visually mine spatial interaction systems using a large domestic air travel dataset. This property can even be used to map data without spatial properties in a 2D or 3D space. For example, [116] mapped massive textual databases with several hundred dimensions in feature space into a 2D space using SOM (Figure 8).

Apart from visualization, one concern about using SOM for geo-referenced data is that the geographic reference of the data is ignored in mapping. [117] suggested a variation of SOM called Geo-SOM to address this issue by considering spatial dependency. Geo-SOM forces the algorithm to search among the neurons geographically close to the data pattern when seeking the winning unit for a specific data pattern. [118] suggests using a one-dimensional SOM to create a sequence of numbers (cluster indices) that are ordered according to the similarity of attributes within the high-dimensional space.



**Figure 7.** (a) Reprinted from [116], Copyright (2013), with permission from Elsevier. The document map of nearly 7 million patent abstracts (background). Twenty best-matched nodes are marked for the query "laser surgery on the cornea" (blues circles within the enlarged map). (b) Reprinted from [118], Copyright (2017), with permission from Elsevier. One-dimensional SOM with 100 nodes, trained with a three-dimensional data set: the final weight vectors of the trained SOM are color-coded based on their index numbers from 1 to 100. (c) Reprinted from [118], Copyright (2017), with permission from Elsevier. Contextual map of London based on the average working hours, the average distance traveled to work, and the average age within each ward.

Each spatial point can then be labeled with its associated index number and represented in a choropleth map. The spatial pattern of these number sequences, called contextual numbers, summarizes the variations of geographic locations in high-dimensional space in a single contextual map (Figure 8b). Like other dimensionality reduction methods in section 4.1.3, such a feature (with the advantage of being a nonlinear dimensionality reduction method) can be used as input to machine learning. For a complete review of SOM applications in GIS, see [115].

#### 4.2.4. Radial Basis Function Networks

Radial basis function (RBF) networks have an input layer, a hidden layer, and an output layer. Instead of a linear relationship between the input vector and the neurons in the hidden layer followed by a nonlinear activation function, the weighted norm (distance) of the input vector and the neurons are calculated in RBF networks with a radially symmetric activation function, which is usually Gaussian. [119] compared the RBF network with MLP networks for modeling urban change, and showed that RBF demonstrates higher prediction accuracy. [120] used the RBF networks for spatial interpolation by incorporating a semivariogram model where the neurons in the hidden layer are the center of the observation points. [121] demonstrates how a hybrid MLP-RBFN network can improve the spatial interpolation results, where MLP and RBF collaborate to fit surfaces of different types.

#### 4.2.5. Adaptive Resonance Theory Networks

Adaptive resonance theory (ART) based networks are a family of neural networks that have been used in spatial interaction flows, crop classification, and land-use change applications [122–124]. ART-based networks are supervised, self-organizing, and self-stabilizing neural networks that can learn fast in nonstationary environments [125]. Fuzzy ARTMAP, which couples ART-based networks with fuzzy logic, is the most famous ART-based network [126]. It includes two input modules:  $Art_a$  and  $Art_b$ , each

with two layers connected by a map field module.  $Art_a$  matches the input vector to the most similar neurons in its second layer. If the vector is not similar to the current neurons in memory, a new neuron is created. This property enables the  $Art_a$  neural network to adaptively change the topology of the network and adds new experiences to the memory.  $Art_b$ , which maintains the class labels, is connected to  $Art_a$  through a map module. However, the Fuzzy ARTMAP depends on the quality of training data and sensitivity to noise and outliers that may be treated as novel patterns. [122] use the adaptive nature of the Fuzzy ARTMAP in forming its topology to account for spatial heterogeneity in land-use change. Their proposed ART-P-MAP considers the land-use change as a regression problem instead of a classification problem and uses the density of training observations as a confidence measure for prediction through a Bayesian decision approach. This approach increases the generalizability of the model and avoids the problem of adding new neurons due to noise and outliers.

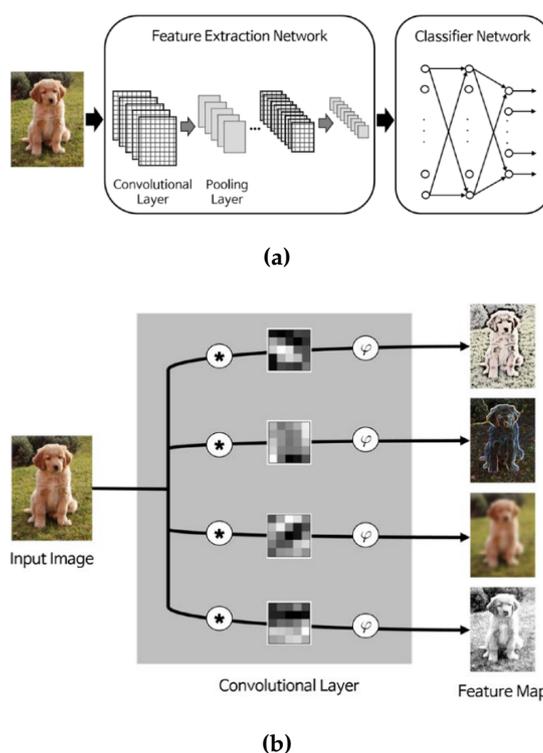
#### 4.2.6. Deep Convolutional Neural Networks

In the past few years, deep neural networks (DNNs) have been proven more promising to process data in their raw form than conventional ML methods. DNNs are usually composed of several nonlinear but simple modules that represent data at different levels. Starting with raw data, each module transforms the representation at one level into a representation at a higher (more abstract) level. In the process and using the backpropagation algorithm, the machine can learn very complex functions [127].

DNNs can be classified into three categories, namely convolutional neural networks (CNNs), generative neural networks (GNNs), and recurrent neural networks (RNNs). Here we discuss several main DNNs that can consider spatial properties of data in their architecture or have been used to solve problems in spatial domains starting with CNNs. RNNs are used primarily to learn from sequence data and will be discussed in section 5.

CNNs are the most popular and well-established form of DNNs to process and analyze images. They include convolutional layers and pooling layers as the two main types Figure 9a. Convolutional layers work based on convolution of a sliding window (filter) with pixel values across the image Figure 9b. The filter weights are determined automatically through the learning process by the network. This is the main advantage of deep CNNs in comparison to conventional ML methods where user defined filter weights are needed. Such feature of CNNs has been used on a limited basis to automatically define W matrix weights in other spatial applications [24]. Pooling layers aggregate neighboring pixels into a single pixel, reducing the image's overall dimensions [128]. As the number of convolutional and pooling layers increases, the network becomes deeper and can extract more abstract features. A classification network finally follows the processing of these layers. In the past few years, several prominent CNN-based DNNs have also been introduced for semantic segmentation (labeling pixels as opposed to image patches) using deep residual networks with less complexity and depthwise separable convolutions, which reduces computation time and the number of network parameters significantly [129–133].

Two main differences exist between deep convolutional networks and conventional ML methods [128]. First, CNN makes the manual spatial feature extraction described in section 4.1.2 an automatic part of the learning algorithm. That is, weights of the convolution filters are not known a priori but are instead determined through the training process, which in some ways resembles the connection weights in ordinary neural networks. These convolution filters partially account for spatial dependence properties by considering the neighborhood pixels, but their size and number within each layer must still be determined in a hyperparameter optimization process. These regularly shaped filters cannot account for variation of spatial dependence in different directions. Second, pooling operations in CNNs allow them to consider hierarchical features in training. The number and architecture of convolutional and pooling layers is arbitrary to



**Figure 8.** Reprinted by permission from [128], Springer Nature Customer Service Centre GmbH: Apress by Kim, P, Copyright (2017). (a) Basic structure of a CNN. (b) convolutional layer.

some degree and can cause unnecessary complexity of model. Thus, one should always optimize these hyperparameters.

The input images have regular shape in CNNs, while actual objects and regions might be irregularly shaped in the real world. An example of that is in urban land use classification. A regular shape input image may ignore the real boundaries of objects [1]. attempts to address this problem by proposing an object-based convolutional neural network (OCNN). This approach involves a two-step process, where image segmentation is initially applied to partition the image into two object categories of linear shaped (e.g., roads) and general shaped (e.g. buildings) objects with homogeneous spectral and spatial properties using the mean-shift method [134]. Two separate CNNs are trained on these categories of objects with different input image sizes. For linear objects, a small input image size and, for general objects, a large input image size perform better. While this approach outperforms regular deep image segmentation networks, it has some limitations. First, it depends to the accuracy and robustness of segmentation process. The extracted partitions do not necessarily represent the actual boundaries of objects. Second, the segmentation itself is a time consuming operation and parameters of segmentation are still needed to be defined manually by the user.

#### 4.2.7. Deep Graph Neural Networks

CNNs have been applied quite extensively for image classification and segmentation. However, many problems (e.g., social and biological networks) can usually not be represented in grid format, making it challenging to apply convolutions. Thus, attempts have been made to extend neural networks to phenomena that are best portrayed with graph structure. Graph neural networks (CNN) were first introduced by [135]. Recently growing attempts have been made to generalize convolution to graphs, which can be categorized into spectral and non-spectral approaches [136]. Spectral methods create a spectral representation for the graph and apply convolution through the graph Fourier Transform [137,138]. The challenge with these types of graphs is, if

the structure of the graph changes, the trained model of the previous structure cannot directly be applied to a graph with new structure. Non-spectral methods directly use convolutions on close neighbors in the graph [139]. These approaches are relatively new and have shown impressive performance in many applications, such as disease spread forecasting [140], traffic analysis [141], medical diagnosis and analysis [142], and natural language processing [143]. A survey of deep learning methods for graphs can be found at [144]. The application of GNNs has yet to be explored in spatial domains especially for non-grid-based spatial data such as social networks.

#### 4.2.8. Deep Generative Networks

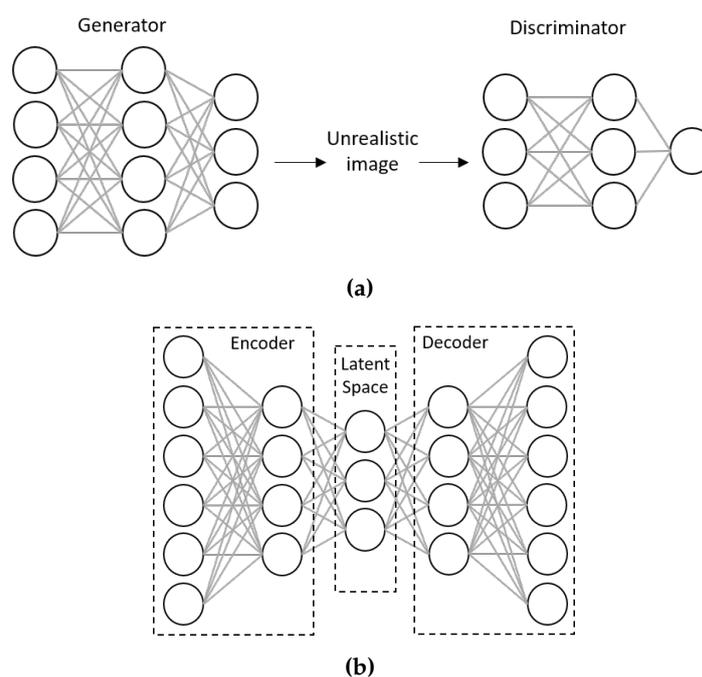
Models in ML can also be categorized as discriminative or generative [145]. On the one hand, we have discriminative models that are usually used for classification and are called classifiers. They return,  $P(Y|X)$ , the probability of a sample to belong to class  $Y$  given the feature attributes  $X$ . Generative models, on the other hand, attempt to generate realistic features of a class of objects, given the distribution of the class  $P(X|Y)$ . Thus, a new set of features is generated using the distribution of a specific class. In the past few years, generative neural network models have attracted considerable attention. Among several types of generative models, variational autoencoders (VAE) (Figure ??) and generative adversarial neural networks (GAN) (Figure ??) have found applications for spatial prediction [146,147].

GANs were first suggested by [146]. They form a type of generative neural networks composed of two sub-networks: a discriminator and a generator. These two networks compete with one another to learn through the training. On the one hand, the generator attempts to create realistic sets of features  $X$  from data class distribution  $Y$  by adding random noise  $Z$  (Figure ??). Random noise  $Z$  is added to make sure each time a new instance  $\hat{X}$  of class  $Y$  is generated. The generator attempts to confuse the discriminator by introducing this new feature set  $\hat{X}$  as a new instance of the class. The discriminator, on the other hand, is a binary classifier that compares  $\hat{X}$  with actual instances of the class and recognizes whether it is real or fake. As training continues, both networks compete until the discriminator fails to recognize the unreal instance from the actual sample. VAE works based on an encoder-decoder architecture, which has the same number of units in the input and output layers [147]. For image classification, for example, VAE creates a new representation of a labeled input image into a space of lower dimensionality; it creates a distribution of the object class in the latent space, and reconstructs a sample image from the distribution. By comparing the actual and reconstructed images, the network can learn through the training.

Many spatial phenomena are heterogeneous and nonlinear, rendering conventional data analytics methods less effective. Generative networks have been applied successfully for DEM spatial interpolation [148], spatiotemporal imputation of aerosol when a substantial amount of data is missing [149], and predicting regional desirability with VGI data [17]. However, CNN-based methods are not appropriate when a large amount of data is missing since they require complete images or images with limited random missing values for training [149]. The other application of generative models is for data augmentation, especially when the size of the training data set is small or the class balance is uneven. Data augmentation works by slightly manipulating the training data to generate new training samples [150].

#### 4.2.9. Learning With DNNs

Apart from discussing the specific capability of the various types of networks, it is also necessary to look at the ways existing DNNs can be used. One can use pre-trained networks, adapt a pre-trained network, or train a network with new data from scratch [151]. In the first method, pre-trained models are used as feature extractors for classification problems on a new data set. This is similar to the methods in section 4.1.2, where new features are created and added to the observation matrix. These methods



**Figure 9.** The basic structure of two generative neural network models. (a) Variational autoencoders (VAE). (b) generative adversarial neural networks (GAN).

have been shown to be effective at classifying remote sensing and photogrammetric imagery [152–154].

Alternatively, one can fine-tune a pre-trained network on a small set of new observations that are sometimes in very different domains or topics. For example, a DNN trained on a specific type of infectious disease (e.g., coronaviruses) to predict the weekly number of cases may be fine-tuned to predict another type of infectious disease (e.g., influenza) with only a small set of data for influenza cases. Another popular example is the trained network on close-range images in urban areas that can be fine-tuned to classify remote sensing images.

The success of pre-trained networks comes from the capability of models to generalize and specifically from the fact that, even though data sets may be from different settings or domains, they still exhibit fundamental standard features that machines can learn from in one setting and apply to another. For example, the geometric boundary of objects in both sets of images mentioned above is composed of corners and horizontal, vertical, and diagonal lines. The way this is usually conducted is that coefficients for several layers in the trained network (usually the layers with low-level features) are frozen, and only the coefficients of the rest of the network are fine-tuned with the new observations. Compared to the first approach, fine-tuning must provide better results since features are oriented to the new data set [155]. Finally, the last method is to train the network on the new data set from scratch. However, this method is subject to overfitting if the training data set is small [151].

#### 4.2.10. Spatial Accuracy

Measures of the accuracy of spatial data are essential to be considered in the objective function in machine learning for a number of geospatial applications. For example, [156] showed how the lack of a spatial accuracy measure could influence the evaluation of location prediction performance for birds nests based on several habitat and environmental factors. When the variables are rasterized, one typical way to evaluate prediction accuracy is to measure the similarity between the predicted and actual maps. In the example being discussed, a cell is labelled as either a nest or no-nest (a binary classification task). Then, the following objective function of learning performance can be

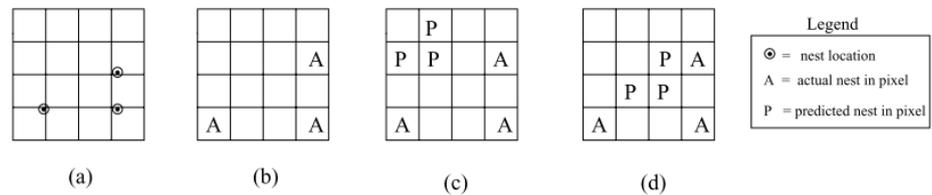
used for this purpose, where a measure is devised to calculate the classification accuracy from the confusion matrix:

$$F = \text{Similarity}(\text{classification accuracy}) \quad (3)$$

Figure 11a and Figure 11b show the location of sample nests and their rasterized format. Figure 11c and Figure ?? represent the predicted locations in two different iterations during learning. Objective function (3) returns the same similarity value for c and d, while d is more accurate spatially. Thus, the authors suggest adding a term to the objective function to measure spatial accuracy, which could be the average distance to the nearest prediction. Therefore, we can rewrite the objective function as follows:

$$F = \text{Similarity}((1 - \alpha) \times \text{classification accuracy} + \alpha \times \text{spatial accuracy}) \quad (4)$$

Where  $\alpha$  is a regularizer parameter that is fine-tuned during the training.



**Figure 10.** Reprinted with permission from [156] (Copyright ©2001 Society for Industrial and Applied Mathematics. All rights reserved.) Spatial accuracy. (a) Actual location of nests, (b) Rasterized location of nests, (c) Predicted locations 1, (d) Predicted locations 2.

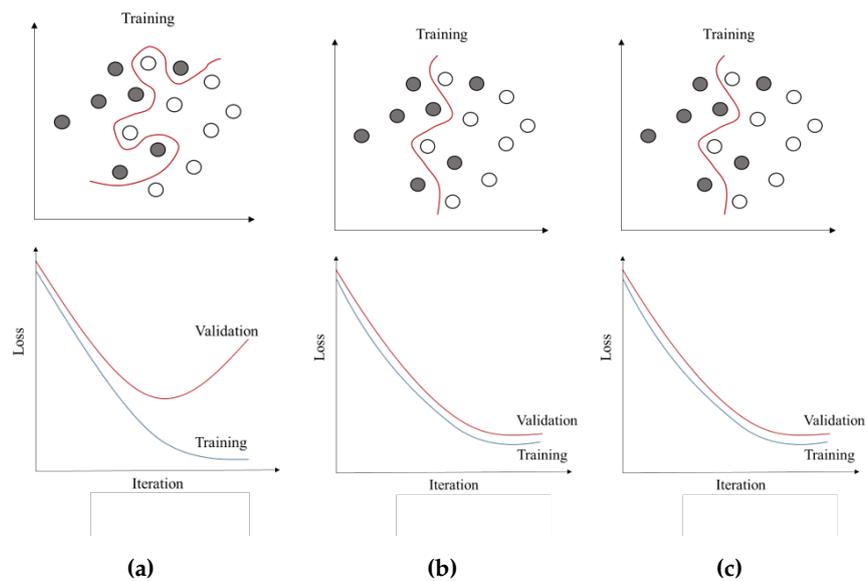
#### 4.2.11. Generalization

Generalization is the capability to generalize a trained model based on a set of data to future datasets. A dataset is usually divided into three mutually exclusive sets, namely a training set, a validation set, and a test set. We fit a model to our training dataset at each iteration and compute an objective function to measure learning performance. Then, we use the validation dataset to evaluate the model's capability to fit another dataset for generalization. This process is repeated until we have a model that fits best on both datasets. The critical point is that the validation dataset is still being used in learning, and we need to ensure that our model is being tested on a completely different dataset. The test dataset is used for the latter purpose, after we select the best model from the previous step. More complex validation methods such as k-fold cross-validation can be used to guarantee the best setting of train, test, and validation dataset [157].

If we fail to reduce the validation loss while the training loss decreases, we will have the so-called overfitting problem (Figure ??). On the other hand, if the model is very simple, both training and validation loss are close but decrease significantly (Figure ??). This means we need to control the complexity of the model to reach the trade-off between the loss value in training and the validation data set. Regularization is the method that is usually used for this purpose, which works by penalizing the weights in the loss function. We can add a regularization term to the objective function as follows:

$$\text{Objective function} : \text{Minimize}(\text{Loss} + \lambda \times \text{reg}) \text{ or, } \text{Maximize}(\text{Similarity} + \lambda \times \text{reg}) \quad (5)$$

Two famous regularization terms are  $L_1$  (sum of the absolute value of the weights) and  $L_2$  (some of the square values of the weights) norms. The former tosses out some of the parameters by imposing zero value to their weights, while the latter assures the weights stay close to zero. Lambda is the parameter that we need to tune to determine the best complexity level for our model and is usually domain-dependent. A small



**Figure 11.** Model complexity and regularization. (a) Overfitting: small Lambda value. (b) Well-fitted: lambda value is tuned. (c) Underfitting: large Lambda value.

lambda value means more weight to our training dataset (Figure ??), and a large value means we are selecting a very simple model (Figure ??). We always need to fine-tune the  $\lambda$  value for the dataset to have a suitable generalization model (Figure ??). Regularization is essential when working with small datasets that are more prone to overfitting and when we have many features that impose computational complexity and noise, which are very common when working with spatial data.

As mentioned above and in section 4.2.8, the  $\alpha$  and  $\lambda$  values are domain-specific, which means they need to be optimized at training time. These parameters are called hyperparameters, which are different from standard parameters so that the latter is determined based on the data set. At the same time, the former is dependent on the properties of the model [158]. Therefore, failure to choose the optimized values for these hyperparameters can cause overfitting. However, hyperparameters are not limited to regularization parameters, and their number is different depending on the machine learning method and problem formulation. For example, random forests have three primary hyperparameters: ensemble size, the maximum size of the individual trees, and the number of randomly selected variables at each node. However, in neural networks, the architecture, learning algorithm, number of training iterations, learning rate, and momentum must be set [159]. In addition, new hyperparameters may be added for spatial applications, such as the spatial accuracy, the initial size and scale of input images, size of the filters, and appropriate sample size. A detailed discussion of hyperparameter optimization in spatial science has been discussed [158].

## 5. Spatiotemporal Learning

A significant amount of attention has been devoted to spatiotemporal learning in the past few years with the availability of technology to collect data at a much higher frequency frequently, if not continuously. Among machine learning methods, neural networks and SVMs have been used often for space-time learning. Similar to the spatial properties discussed in this paper, spatiotemporal dependence and non-stationarity may also exist in data with spatiotemporal dimensions [160]. The number of parameters in spatiotemporal ML may become very large, which can make learning impossible if the model cannot capture the underlying spatiotemporal structure well [161].

Geographically and temporally weighted regression models have already been developed for geospatial applications. Still, the challenges related to expressing complex and nonlinear space-time proximity and optimal weights for kernels remain unanswered

in these methods. [24] proposed a spatiotemporal proximity neural network (STPNN) that constructs the nonstationary weight matrix instead of using fixed and conventional methods to address the complex nonlinear interactions between time and space. [22] used a multi-stage approach to address spatial heterogeneity and dependence for space-time prediction. Authors used GeoSoM to divide space and time into homogeneous regions. In the second step of the latter model, a space-time lag within each cluster was estimated to capture the space-time dependence structure among the space-time series. Finally, a feedback recurrent neural network predicts values on each cluster locally. Although such techniques have high performance, they are usually multistep, and the computational cost is high. Furthermore, complexities related to anisotropy are not modeled.

Convolutional recurrent neural networks and especially convolutional long short term memory (LSTM) networks can be applied extensively in spatiotemporal learning of grid data [162]. LSTM is a type of recurrent neural network with the capability to memorizing temporal dependencies in data. A combination of this feature with the power of CNNs to learn the hierarchical spatial features can provide an automatic single-step ML model to account for space-time dependence. A recent survey of machine learning methods of spatiotemporal sequences is available elsewhere for interested readers [161].

## 6. Discussion and Concluding Remarks

We conducted a state-of-the-art survey of literature where ML cross-pollinates spatial domains in which data exhibit distinctive properties such as spatial dependence, spatial heterogeneity, and scale. We identified two broad approaches in this body of literature, which are respectively motivated by the two components of a spatial ML learning system, namely the spatial observation matrix and the learning algorithm. The former explicitly handles the spatial properties of data before the process of learning begins. In other words, no change in the learning algorithm is implemented following this step. It is now well recognized that considering spatial properties in sampling strategies and in addressing missing data is necessary in any spatial ML application. In addition to these matters, creating new spatial features was discussed as one of the main approaches to augment the observation matrix with new spatial properties of data. To date, a large body of literature in ML of spatially explicit data has resorted to spatial features mainly because the idea comes naturally, because extensive research in geographic information science has focused on these matters over the past two decades, and because this approach permits using existing ML algorithms without further modifications. Many of these methods have successfully been used for a variety of applications, ranging from point cloud classification to trajectory analysis and pattern recognition in satellite imagery.

We also discussed how spatial properties can be handled explicitly in the other component of ML, namely the learning algorithm, an approach that has only recently started to be explored. Here, spatial properties are addressed in the learning algorithm representation or objective function rather than at the level of the observation matrix. When dealing with learning algorithms in spatial domains, we argued for focused attention to spatial hyperparameter optimization and spatial accuracy. Different learning algorithms require various numbers of hyperparameters to be optimized, with deep learning methods usually having the largest number. When it comes to accuracy, spatial accuracy is often ignored, while evidence shows that it can significantly influence the results and the generalizability of the model, and from there, degrade the predictive power of the ML model. New measures for spatial accuracy may be needed to alleviate these issues. Space-time learning has recently also become a focus of considerable attention, both in identifying technical challenges and in advancing modeling solutions, as processes intermingle in space and time. With the proliferation of panel and other

space-time data and the focused interest for process-based knowledge triggered by the COVID-19 outbreak and pandemic, this area has emerged as a priority research area.

Our literature review shows that progress in the learning algorithm component of ML is still in early stages compared to advances made with enhancing the spatial observation matrix, and there is a lot more room to develop and apply some of these methods in different spatial domains. Here, the main takeaways are as follows.

- CNNs can be used to automatically estimate the spatial weight matrix, which is usually unknown and needs to be defined by the user to reflect spatial data properties in many spatial problems. Advances may be anticipated in several areas;
- Deep neural networks with convolutional layers have been shown to automatically extract patterns from multiple scales and hierarchies. However, they have so far mainly been used to recognize patterns in raster data sets, and use cases in a broader range of domains of applications are called for;
- Graph-based deep learning methods provide a new opportunity to apply CNN-based deep learning to problems with graph structure (e.g., social networks) or when the geographic units are irregularly shaped (e.g., census data);
- Further studies for learning in spatio-temporal domains will need to be undertaken as well. Deep neural networks based on a combination of LSTM and CNNs introduce simultaneous learning across space, time, scales, and hierarchies. When augmented with reinforcement learning to add feedback within systems, which is the case in many spatial, social and environmental applications, they can realize the dream of a single universal ML method [127].

However, deep neural networks have their own limitations and pitfalls. First, they need a large amount of training data. Second, they have a large number of parameters and hyperparameters, which make them computationally expensive. Third, there is no limit to increasing the complexity of the DNNs, because of the arbitrary nature of architecture design. Thus, there are always potential pitfalls for researchers to develop less reproducible and unnecessarily complex models.

In addition to the above suggestions for future research, which both extend past research practices and leverage some of the proposed methods discussed in this paper to apply ML for spatial data analytics, a long-term line of research in this area may allign with a more fundamental recasting of concepts related to the definition of space. Properties of spatial data that we discussed in this paper are the byproducts of the current conceptual definition of space as a fixed “container”. This container is independent of the objects and events that exist and occur inside it. With the current approach, space is defined in a reference coordinate system tied to an origin with a scale measure, where the location of entities is determined with respect to this origin. This is how all of the current GIS and Remote Sensing software and tools conceptualize space, and by extension, so does spatial ML.

While this absolute view of space has been found practical thus far, it has its own limitations. For example, the spatial weight matrix, designed to involve the neighborhood relationships between geographic entities, may be an inadequate representation of space. It is independent of time, usually assumes isotropy in neighborhood definition, and can only partially account for spatial dependence. Also, it cannot capture the complete complexity of spatial problems when multiple variables are present. Another example of that is spatial lags, which are limited since they change for different variables. If one wants to include them in learning, the final model may end up with a very large number of variables. At the same time, the user still needs to discover how many spatial lags and for which variables to use them. Such methods cannot address more complex problems such as UGCoP and MAUP, making it difficult to predetermine scales and appropriate units of observation. The contextual influences naturally change across granularities, which makes the problem even more complex.

Contrary to the absolute view, a relative view of space assumes that space is a construct of the events that take place within it and is not completely separate from them

[40]. From this point of view, Euclidian coordinates are not favored over other semantic attributes, and they collectively construct the spatial space and structure. Current tools and software require a reconceptualization of the space to represent this relative view as no easy and robust tuning of existing paradigms can be envisioned. From this point of view, we see attractive and thus far almost untouched areas that primarily land in the second component of Figure 2, the learning algorithm. For example, graph-based deep learning has remarkable potential for future research within this broader view. Graphs can provide a good representation of relative space mainly because they can explicitly model the relations and are not required to be in a common coordinate reference frame. Current GIS tools and software need all layers to be in the same reference coordinate system and they perform quite poorly at handling networks of social interactions. No limit exists in the definition of neighbors in graphs, and they can dynamically change. Nested and hierarchical networks with different geometric references and topology are easy to link (e.g., in bipartite networks) and manage. MAUP disappears because objects and phenomena need not to be aggregated in partitioned areal units anymore. Nested and hierarchical graphs may also be a solution to the UGCoP. One can actually delineate the boundaries of environmental risk exposures for individuals who visit different locations (home, school, etc.) during the day within graphs [61].

Graph-based deep learning can become a gateway for machines to learn in social networks, where relations are the building blocks of everything. It is worth noting that the phrase “social networks” is not equivalent to online social media platforms, but as a broader term, any network with human, biological, and natural social interactions. One attractive example is the weighted stochastic block model [163,164], which is a powerful approach to detecting communities in social networks. The spatial organization of communities in a weighted stochastic block model can reveal new relationships in data. Another example is SOMs that organize clusters in space and have widely been studied in the geography community.

The geometric attributes such as coordinates can be treated like other attributes within a relative space, so we do not expect ML methods to change. At the same time, the concept of relative space presents fundamental challenges and questions that need to be answered in future research, such as: Are graphs an adequate data model to represent such a space, or do we need new data models? How can one acquire the ground truth required for most ML methods in a space that is elastic and changing dynamically? What can machines learn about a space that is a construct of its components?

**Author Contributions:** “Conceptualization, B.N. and J.C.T.; methodology, B.N. and J.C.T.; formal analysis, B.N. and J.C.T.; investigation, B.N. and J.C.T.; resources, B.N. and J.C.T.; writing—original draft preparation, B.N. and J.C.T.; writing—review and editing, B.N. and J.C.T.; visualization, B.N.; supervision, J.C.T.; project administration, J.C.T.; All authors have read and agreed to the published version of the manuscript.”

**Funding:** Please add: “This research received no external funding”.

**Conflicts of Interest:** “The authors declare no conflict of interest.”

## References

1. Zhang, C.; Sargent, I.; Pan, X.; Li, H.; Gardiner, A.; Hare, J.; Atkinson, P.M. An object-based convolutional neural network (OCNN) for urban land use classification. *Remote sensing of environment* **2018**, *216*, 57–70.
2. Zhang, C.; Sargent, I.; Pan, X.; Li, H.; Gardiner, A.; Hare, J.; Atkinson, P.M. Joint Deep Learning for land cover and land use classification. *Remote sensing of environment* **2019**, *221*, 173–187.
3. Law, S.; Seresinhe, C.I.; Shen, Y.; Gutierrez-Roig, M. Street-Frontage-Net: urban image classification using deep convolutional neural networks. *International Journal of Geographical Information Science* **2020**, *34*, 681–707.
4. Srivastava, S.; Vargas Munoz, J.E.; Lobry, S.; Tuia, D. Fine-grained landuse characterization using ground-based pictures: a deep learning solution based on globally available data. *International Journal of Geographical Information Science* **2020**, *34*, 1117–1136.
5. Hagenauer, J.; Omrani, H.; Helbich, M. Assessing the performance of 38 machine learning models: the case of land consumption rates in Bavaria, Germany. *International Journal of Geographical Information Science* **2019**, *33*, 1399–1419.

6. Guan, Q.; Wang, L.; Clarke, K.C. An artificial-neural-network-based, constrained CA model for simulating urban growth. *Cartography and Geographic Information Science* **2005**, *32*, 369–380.
7. Reades, J.; De Souza, J.; Hubbard, P. Understanding urban gentrification through machine learning. *Urban Studies* **2019**, *56*, 922–942.
8. Resch, B.; Usländer, F.; Havas, C. Combining machine-learning topic models and spatiotemporal analysis of social media data for disaster footprint and damage assessment. *Cartography and Geographic Information Science* **2018**, *45*, 362–376.
9. Masjedi, A.; Crawford, M.M. PREDICTION OF SORGHUM BIOMASS USING TIME SERIES UAV-BASED HYPERSPECTRAL AND LIDAR DATA. IGARSS 2020-2020 IEEE International Geoscience and Remote Sensing Symposium. IEEE, 2020, pp. 3912–3915.
10. Adhikari, B.; Xu, X.; Ramakrishnan, N.; Prakash, B.A. Epideep: Exploiting embeddings for epidemic forecasting. Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 2019, pp. 577–586.
11. Effati, M.; Thill, J.C.; Shabani, S. Geospatial and machine learning techniques for wicked social science problems: analysis of crash severity on a regional highway corridor. *Journal of Geographical Systems* **2015**, *17*, 107–135.
12. Skupin, A.; Hagelman, R. Visualizing demographic trajectories with self-organizing maps. *Geoinformatica* **2005**, *9*, 159–179.
13. Steiniger, S.; Taillandier, P.; Weibel, R. Utilising urban context recognition and machine learning to improve the generalisation of buildings. *International Journal of Geographical Information Science* **2010**, *24*, 253–282.
14. Cunha, E.; Martins, B. Using one-class classifiers and multiple kernel learning for defining imprecise geographic regions. *International Journal of Geographical Information Science* **2014**, *28*, 2220–2241.
15. Chegoonian, A.; Mokhtarzade, M.; Valadan Zoj, M. A comprehensive evaluation of classification algorithms for coral reef habitat mapping: challenges related to quantity, quality, and impurity of training samples. *International Journal of Remote Sensing* **2017**, *38*, 4224–4243.
16. Lin, Y.; Kang, M.; Wu, Y.; Du, Q.; Liu, T. A deep learning architecture for semantic address matching. *International Journal of Geographical Information Science* **2020**, *34*, 559–576.
17. Shi, W.; Liu, Z.; An, Z.; Chen, P. RegNet: a neural network model for predicting regional desirability with VGI data. *International Journal of Geographical Information Science* **2021**, *35*, 175–192.
18. Yang, C.; Gidófalvi, G. Detecting regional dominant movement patterns in trajectory data with a convolutional neural network. *International Journal of Geographical Information Science* **2020**, *34*, 996–1021.
19. Zhao, R.; Pang, M.; Wang, J. Classifying airborne LiDAR point clouds via deep features learned by a multi-scale convolutional neural network. *International journal of geographical information science* **2018**, *32*, 960–979.
20. Yan, J.; Thill, J.C. Visual data mining in spatial interaction analysis with self-organizing maps. *Environment and Planning B: Planning and Design* **2009**, *36*, 466–486.
21. Rigol, J.P.; Jarvis, C.H.; Stuart, N. Artificial neural networks as a tool for spatial interpolation. *International Journal of Geographical Information Science* **2001**, *15*, 323–343.
22. Deng, M.; Yang, W.; Liu, Q. Geographically weighted extreme learning machine: a method for space–time prediction. *Geographical Analysis* **2017**, *49*, 433–450.
23. Deng, M.; Yang, W.; Liu, Q.; Jin, R.; Xu, F.; Zhang, Y. Heterogeneous space–time artificial neural networks for space–time series prediction. *Transactions in GIS* **2018**, *22*, 183–201.
24. Wu, S.; Wang, Z.; Du, Z.; Huang, B.; Zhang, F.; Liu, R. Geographically and temporally neural network weighted regression for modeling spatiotemporal non-stationary relationships. *International Journal of Geographical Information Science* **2021**, *35*, 582–608.
25. Koperski, K.; Adhikary, J.; Han, J. Spatial data mining: progress and challenges survey paper. Proc. ACM SIGMOD Workshop on Research Issues on Data Mining and Knowledge Discovery, Montreal, Canada. Citeseer, 1996, pp. 1–10.
26. Mennis, J.; Guo, D. Spatial data mining and geographic knowledge discovery—An introduction. *Computers, Environment and Urban Systems* **2009**, *33*, 403–408.
27. Miller, H.J.; Han, J. *Geographic data mining and knowledge discovery*; CRC press, 2009.
28. Jiang, Z. A survey on spatial prediction methods. *IEEE Transactions on Knowledge and Data Engineering* **2018**, *31*, 1645–1664.
29. Gopal, S. Artificial neural networks in geospatial analysis. *International Encyclopedia of Geography: People, the Earth, Environment and Technology: People, the Earth, Environment and Technology* **2016**, pp. 1–7.
30. Kanevski, M. *Machine learning for spatial environmental data: theory, applications, and software*; EPFL press, 2009.
31. Yang, L.; MacEachren, A.M.; Mitra, P.; Onorati, T. Visually-enabled active deep learning for (geo) text and image classification: a review. *ISPRS International Journal of Geo-Information* **2018**, *7*, 65.
32. Mitchell, T.M.; others. *Machine learning*; Vol. 1, McGraw-Hill Science, 1997.
33. Friedman, J.; Hastie, T.; Tibshirani, R.; others. *The elements of statistical learning*; Vol. 1, Springer series in statistics New York, 2001.
34. Ghahramani, Z. Unsupervised learning. Summer School on Machine Learning. Springer, 2003, pp. 72–112.
35. Zhu, X.; Goldberg, A.B. Introduction to semi-supervised learning. *Synthesis lectures on artificial intelligence and machine learning* **2009**, *3*, 1–130.
36. Settles, B. *Active learning literature survey*; Vol. 1, University of Wisconsin-Madison Department of Computer Sciences, 2009.
37. Haining, R. The special nature of spatial data. *The SAGE Handbook of Spatial Analysis*. Los Angeles: SAGE Publications **2009**, pp. 5–23.

38. Anselin, L. What is Special About Spatial Data? Alternative Perspectives on Spatial Data Analysis (89-4). *eScholarship University of California Santa Barbara* **1989**.
39. Getis, A. Spatial dependence and heterogeneity and proximal databases. *Spatial analysis and GIS* **1994**, pp. 105–120.
40. Thill, J.C. Is spatial really that special? A tale of spaces. In *Information Fusion and Geographic Information Systems*; Springer, 2011; pp. 3–12.
41. Tobler, W.R. A computer movie simulating urban growth in the Detroit region. *Economic geography* **1970**, *46*, 234–240.
42. Getis, A. Spatial autocorrelation. In *Handbook of applied spatial analysis*; Springer, 2010; pp. 255–278.
43. Griffith, D.A. Spatial autocorrelation. *A Primer (Washington, DC, Association of American Geographers)* **1987**.
44. Moran, P.A. Notes on continuous stochastic phenomena. *Biometrika* **1950**, *37*, 17–23.
45. Geary, R.C. The contiguity ratio and statistical mapping. *The incorporated statistician* **1954**, *5*, 115–146.
46. Cressie, N. *Statistics for spatial data*; John Wiley & Sons, 2015.
47. Delmelle, E. Spatial sampling. *The SAGE handbook of spatial analysis* **2009**, 183, 206.
48. Dao, T.H.D.; Thill, J.C. The SpatialARMED framework: handling complex spatial components in spatial association rule mining. *Geographical Analysis* **2016**, *48*, 248–274.
49. Dale, M.R.; Fortin, M.J. *Spatial analysis: a guide for ecologists*; Cambridge University Press, 2014.
50. Thill, J.C. Research on urban and regional systems: Contributions from GIS, spatial analysis, and location modeling. *Innovations in Urban and Regional Systems* **2020**, pp. 3–20.
51. Murwira, A.; Skidmore, A.K. The response of elephants to the spatial heterogeneity of vegetation in a Southern African agricultural landscape. *Landscape ecology* **2005**, *20*, 217–234.
52. Webster, R. Is soil variation random? *Geoderma* **2000**, *97*, 149–163.
53. Hu, Y.; Li, W.; Wright, D.; Aydin, O.; Wilson, D.; Maher, O.; Raad, M. Artificial intelligence approaches. *arXiv preprint arXiv:1908.10345* **2019**.
54. Lam, N. FC-21-Resolution. *University Consortium for Geographic Information Science GIS and T Body of Knowledge*.
55. Shekhar, S.; Gandhi, V.; Zhang, P.; Vatsavai, R.R.; Fotheringham, A.; Rogerson, P. Availability of spatial data mining techniques. *The SAGE handbook of spatial analysis* **2009**, pp. 159–181.
56. Fotheringham, A.S.; Wong, D.W. The modifiable areal unit problem in multivariate statistical analysis. *Environment and planning A* **1991**, *23*, 1025–1044.
57. Openshaw, S. A million or so correlation coefficients, three experiments on the modifiable areal unit problem. *Statistical applications in the spatial science* **1979**, pp. 127–144.
58. Arbia, G. *Spatial data configuration in statistical analysis of regional economic and related problems*; Vol. 14, Springer Science & Business Media, 2012.
59. Batty, M.; Sikdar, P. Spatial aggregation in gravity models. 1. An information-theoretic framework. *Environment and Planning A* **1982**, *14*, 377–405.
60. Xiao, J. Spatial Aggregation Entropy: A Heterogeneity and Uncertainty Metric of Spatial Aggregation. *Annals of the American Association of Geographers* **2021**, *111*, 1236–1252.
61. Kwan, M.P. The uncertain geographic context problem. *Annals of the Association of American Geographers* **2012**, *102*, 958–968.
62. Robinson, W.S. Ecological correlations and the behavior of individuals. *International journal of epidemiology* **2009**, *38*, 337–341.
63. Zeng, W.; Lin, C.; Lin, J.; Jiang, J.; Xia, J.; Turkay, C.; Chen, W. Revisiting the modifiable areal unit problem in deep traffic prediction with visual analytics. *IEEE Transactions on Visualization and Computer Graphics* **2020**, *27*, 839–848.
64. Chawla, S.; Shekhar, S.; Wu, W.; Özesmi, U. Predicting Locations Using Map Similarity (PLUMS): A Framework for Spatial Data Mining. *MDM/KDD*, 2000, pp. 14–24.
65. Chen, J.; Lu, F.; Peng, G. A quantitative approach for delineating principal fairways of ship passages through a strait. *Ocean Engineering* **2015**, *103*, 188–197.
66. Acheson, E.; Volpi, M.; Purves, R.S. Machine learning for cross-gazetteer matching of natural features. *International Journal of Geographical Information Science* **2020**, *34*, 708–734.
67. Santos, R.; Murrieta-Flores, P.; Calado, P.; Martins, B. Toponym matching through deep neural networks. *International Journal of Geographical Information Science* **2018**, *32*, 324–348.
68. Purves, R.; Jones, C. Geographic information retrieval. *SIGSPATIAL Special* **2011**, *3*, 2–4.
69. Yao, X.; Thill, J.C. Spatial queries with qualitative locations in spatial information systems. *Computers, environment and urban systems* **2006**, *30*, 485–502.
70. Guo, Z.; Feng, C.C. Using multi-scale and hierarchical deep convolutional features for 3D semantic classification of TLS point clouds. *International Journal of Geographical Information Science* **2020**, *34*, 661–680.
71. Li, X.; Ma, R.; Zhang, Q.; Li, D.; Liu, S.; He, T.; Zhao, L. Anisotropic characteristic of artificial light at night—Systematic investigation with VIIRS DNB multi-temporal observations. *Remote Sensing of Environment* **2019**, *233*, 111357.
72. Özyeşil, O.; Voroninski, V.; Basri, R.; Singer, A. A survey of structure from motion\*. *Acta Numerica* **2017**, *26*, 305–364.
73. Fuentes-Pacheco, J.; Ruiz-Ascencio, J.; Rendón-Mancha, J.M. Visual simultaneous localization and mapping: a survey. *Artificial intelligence review* **2015**, *43*, 55–81.

74. Jain, S.; Smit, A.; Truong, S.Q.; Nguyen, C.D.; Huynh, M.T.; Jain, M.; Young, V.A.; Ng, A.Y.; Lungren, M.P.; Rajpurkar, P. VisualCheXbert: addressing the discrepancy between radiology report labels and image labels. *Proceedings of the Conference on Health, Inference, and Learning*, 2021, pp. 105–115.
75. Zhang, J.; Liu, J.; Pan, B.; Shi, Z. Domain adaptation based on correlation subspace dynamic distribution alignment for remote sensing image scene classification. *IEEE Transactions on Geoscience and Remote Sensing* **2020**, *58*, 7920–7930.
76. Gahegan, M. Is inductive machine learning just another wild goose (or might it lay the golden egg)? *International Journal of Geographical Information Science* **2003**, *17*, 69–92.
77. Congalton, R.G. A review of assessing the accuracy of classifications of remotely sensed data. *Remote sensing of environment* **1991**, *37*, 35–46.
78. Pontius Jr, R.G.; Millones, M. Death to Kappa: birth of quantity disagreement and allocation disagreement for accuracy assessment. *International Journal of Remote Sensing* **2011**, *32*, 4407–4429.
79. Buda, M.; Maki, A.; Mazurowski, M.A. A systematic study of the class imbalance problem in convolutional neural networks. *Neural Networks* **2018**, *106*, 249–259.
80. Hase, N.; Ito, S.; Kaneko, N.; Sumi, K. Data augmentation for intra-class imbalance with generative adversarial network. Fourteenth International Conference on Quality Control by Artificial Vision. International Society for Optics and Photonics, 2019, Vol. 11172, p. 1117206.
81. Martin, R.; Aler, R.; Valls, J.M.; Galván, I.M. Machine learning techniques for daily solar energy prediction and interpolation using numerical weather models. *Concurrency and Computation: Practice and Experience* **2016**, *28*, 1261–1274.
82. Zanella, L.; Folkard, A.M.; Blackburn, G.A.; Carvalho, L.M. How well does random forest analysis model deforestation and forest fragmentation in the Brazilian Atlantic forest? *Environmental and ecological statistics* **2017**, *24*, 529–549.
83. Meyer, H.; Reudenbach, C.; Wöllauer, S.; Nauss, T. Importance of spatial predictor variable selection in machine learning applications—Moving from data reproduction to spatial prediction. *Ecological Modelling* **2019**, *411*, 108815.
84. Anselin, L.; Arribas-Bel, D. Spatial fixed effects and spatial dependence in a single cross-section. *Papers in Regional Science* **2013**, *92*, 3–17.
85. Sommervoll, Å.; Sommervoll, D.E. Learning from man or machine: Spatial fixed effects in urban econometrics. *Regional Science and Urban Economics* **2019**, *77*, 239–252.
86. Wu, S.S.; Qiu, X.; Usery, E.L.; Wang, L. Using geometrical, textural, and contextual information of land parcels for classification of detailed urban land use. *Annals of the Association of American Geographers* **2009**, *99*, 76–98.
87. Cristóbal, G.; Bescós, J.; Santamaría, J. Image analysis through the Wigner distribution function. *Applied optics* **1989**, *28*, 262–271.
88. Myint, S.W. A robust texture analysis and classification approach for urban land-use and land-cover feature discrimination. *Geocarto international* **2001**, *16*, 29–40.
89. Turner, M. Texture transformation by Gabor function. *Biology Cybernation* **1986**, *55*, 71–82.
90. Zhu, C.; Yang, X. Study of remote sensing image texture analysis and classification using wavelet. *International Journal of Remote Sensing* **1998**, *19*, 3197–3203.
91. Platt, R.V.; Rapoza, L. An evaluation of an object-oriented paradigm for land use/land cover classification. *The Professional Geographer* **2008**, *60*, 87–100.
92. Zhan, Q.; Molenaar, M.; Gorte, B. Urban land use classes with fuzzy membership and classification based on integration of remote sensing and GIS. *International Archives of Photogrammetry and Remote Sensing* **2000**, *33*, 1751–1759.
93. Herold, M.; Liu, X.; Clarke, K.C. Spatial metrics and image texture for mapping urban land use. *Photogrammetric Engineering & Remote Sensing* **2003**, *69*, 991–1001.
94. Zhou, Y.; Tuzel, O. Voxelnet: End-to-end learning for point cloud based 3d object detection. *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4490–4499.
95. Lawin, F.J.; Danelljan, M.; Tosteberg, P.; Bhat, G.; Khan, F.S.; Felsberg, M. Deep projective 3D semantic segmentation. *International Conference on Computer Analysis of Images and Patterns*. Springer, 2017, pp. 95–107.
96. Yu, S.; Jia, S.; Xu, C. Convolutional neural networks for hyperspectral image classification. *Neurocomputing* **2017**, *219*, 88–98.
97. Hagenauer, J.; Helbich, M. Mining urban land-use patterns from volunteered geographic information by means of genetic algorithms and artificial neural networks. *International Journal of Geographical Information Science* **2012**, *26*, 963–982.
98. Comon, P. Independent component analysis, a new concept? *Signal processing* **1994**, *36*, 287–314.
99. Jolliffe, I.T. Principal components in regression analysis. In *Principal component analysis*; Springer, 1986; pp. 129–155.
100. Kohonen, T. *Self-organization and associative memory*; Vol. 8, Springer Science & Business Media, 2012.
101. Harris, P.; Brunson, C.; Charlton, M. Geographically weighted principal components analysis. *International Journal of Geographical Information Science* **2011**, *25*, 1717–1736.
102. Fotheringham, A.S.; Brunson, C.; Charlton, M. *Geographically weighted regression: the analysis of spatially varying relationships*; John Wiley & Sons, 2003.
103. Qu, L.; Li, L.; Zhang, Y.; Hu, J. PPCA-based missing data imputation for traffic flow volume: A systematical approach. *IEEE Transactions on intelligent transportation systems* **2009**, *10*, 512–522.
104. Cressie, N. The origins of kriging. *Mathematical geology* **1990**, *22*, 239–252.
105. Harris, P.; Fotheringham, A.; Crespo, R.; Charlton, M. The use of geographically weighted regression for spatial prediction: an evaluation of models using simulated data sets. *Mathematical Geosciences* **2010**, *42*, 657–680.

106. Tipping, M.E.; Bishop, C.M. Probabilistic principal component analysis. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **1999**, *61*, 611–622.
107. Chen, H. Principal component analysis with missing data and outliers. *Electrical and Computer Engineering Department Rutgers University* **2002**.
108. Li, Y.; Li, Z.; Li, L.; Zhang, Y.; Jin, M. Comparison on PPCA, KPPCA and MPPCA based missing data imputing for traffic flow. In *ICTIS 2013: Improving Multimodal Transportation Systems-Information, Safety, and Integration*; American Society of Civil Engineers, 2013; pp. 1151–1156.
109. Li, L.; Li, Y.; Li, Z. Efficient missing data imputing for traffic flow by considering temporal and spatial dependence. *Transportation research part C: emerging technologies* **2013**, *34*, 108–120.
110. Stojanova, D.; Ceci, M.; Appice, A.; Malerba, D.; Džeroski, S. Global and local spatial autocorrelation in predictive clustering trees. *International Conference on Discovery Science*. Springer, 2011, pp. 307–322.
111. Jiang, Z.; Shekhar, S.; Zhou, X.; Knight, J.; Corcoran, J. Focal-test-based spatial decision tree learning: A summary of results. 2013 IEEE 13th International Conference on Data Mining. IEEE, 2013, pp. 320–329.
112. Anselin, L. Local indicators of spatial association—LISA. *Geographical analysis* **1995**, *27*, 93–115.
113. Vapnik, V. *The nature of statistical learning theory*; Springer science & business media, 2013.
114. Lee, C.H.; Greiner, R.; Schmidt, M. Support vector random fields for spatial classification. *European Conference on Principles of Data Mining and Knowledge Discovery*. Springer, 2005, pp. 121–132.
115. Agarwal, P.; Skupin, A. *Self-organising maps: Applications in geographic information science*; John Wiley & Sons, 2008.
116. Kohonen, T. Essentials of the self-organizing map. *Neural networks* **2013**, *37*, 52–65.
117. Bação, F.; Lobo, V.; Painho, M. The self-organizing map, the Geo-SOM, and relevant variants for geosciences. *Computers & geosciences* **2005**, *31*, 155–163.
118. Moosavi, V. Contextual mapping: Visualization of high-dimensional spatial patterns in a single geo-map. *Computers, Environment and Urban Systems* **2017**, *61*, 1–12.
119. Shafizadeh-Moghadam, H.; Hagenauer, J.; Farajzadeh, M.; Helbich, M. Performance analysis of radial basis function networks and multi-layer perceptron networks in modeling urban change: a case study. *International Journal of Geographical Information Science* **2015**, *29*, 606–623.
120. Lin, G.F.; Chen, L.H. A spatial interpolation method based on radial basis function networks incorporating a semivariogram model. *Journal of Hydrology* **2004**, *288*, 288–298.
121. Yeh, I.C.; Huang, K.C.; Kuo, Y.H. Spatial interpolation using MLP-RBFN hybrid networks. *International Journal of Geographical Information Science* **2013**, *27*, 1884–1901.
122. Gong, Z.; Thill, J.C.; Liu, W. ART-P-MAP neural networks modeling of land-use change: accounting for spatial heterogeneity and uncertainty. *Geographical Analysis* **2015**, *47*, 376–409.
123. Malamiri, H.R.G.; Aliabad, F.A.; Shojaei, S.; Morad, M.; Band, S.S. A study on the use of UAV images to improve the separation accuracy of agricultural land areas. *Computers and Electronics in Agriculture* **2021**, *184*, 106079.
124. Yariyan, P.; Ali Abbaspour, R.; Chehregan, A.; Karami, M.; Cerdà, A. GIS-based seismic vulnerability mapping: a comparison of artificial neural networks hybrid models. *Geocarto International* **2021**, pp. 1–24.
125. Carpenter, G.A.; Grossberg, S.; Reynolds, J.H. ARTMAP: Supervised real-time learning and classification of nonstationary data by a self-organizing neural network. *Neural networks* **1991**, *4*, 565–588.
126. Carpenter, G.A.; Grossberg, S.; Markuzon, N.; Reynolds, J.H.; Rosen, D.B.; others. Fuzzy ARTMAP: A neural network architecture for incremental supervised learning of analog multidimensional maps. *IEEE Transactions on neural networks* **1992**, *3*, 698–713.
127. LeCun, Y.; Bengio, Y.; Hinton, G. Deep learning. *nature* **2015**, *521*, 436–444.
128. Kim, P. Matlab deep learning. *With machine learning, neural networks and artificial intelligence* **2017**, 130.
129. Chen, L.C.; Papandreou, G.; Kokkinos, I.; Murphy, K.; Yuille, A.L. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence* **2017**, *40*, 834–848.
130. Chollet, F. Xception: Deep learning with depthwise separable convolutions. *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1251–1258.
131. He, K.; Zhang, X.; Ren, S.; Sun, J. Deep residual learning for image recognition *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2016.
132. Ronneberger, O.; Fischer, P.; Brox, T. U-net: Convolutional networks for biomedical image segmentation. *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.
133. Simonyan, K.; Zisserman, A. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* **2014**.
134. Comaniciu, D.; Meer, P. Mean shift: A robust approach toward feature space analysis. *IEEE Transactions on pattern analysis and machine intelligence* **2002**, *24*, 603–619.
135. Gori, M.; Monfardini, G.; Scarselli, F. A new model for learning in graph domains. *Proceedings. 2005 IEEE International Joint Conference on Neural Networks*, 2005. IEEE, 2005, Vol. 2, pp. 729–734.
136. Veličković, P.; Cucurull, G.; Casanova, A.; Romero, A.; Lio, P.; Bengio, Y. Graph attention networks. *arXiv preprint arXiv:1710.10903* **2017**.

137. Estrach, J.B.; Zaremba, W.; Szlam, A.; LeCun, Y. Spectral networks and deep locally connected networks on graphs. 2nd International Conference on Learning Representations, ICLR, 2014, Vol. 2014.
138. Kipf, T.N.; Welling, M. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907* **2016**.
139. Hamilton, W.L.; Ying, R.; Leskovec, J. Inductive representation learning on large graphs. Proceedings of the 31st International Conference on Neural Information Processing Systems, 2017, pp. 1025–1035.
140. Wu, D.; Gao, L.; Xiong, X.; Chinazzi, M.; Vespignani, A.; Ma, Y.A.; Yu, R. DeepGLEAM: a hybrid mechanistic and deep learning model for COVID-19 forecasting. *arXiv preprint arXiv:2102.06684* **2021**.
141. Ye, J.; Zhao, J.; Ye, K.; Xu, C. How to build a graph-based deep learning architecture in traffic domain: A survey. *IEEE Transactions on Intelligent Transportation Systems* **2020**.
142. Ahmedt-Aristizabal, D.; Armin, M.A.; Denman, S.; Fookes, C.; Petersson, L. Graph-Based Deep Learning for Medical Diagnosis and Analysis: Past, Present and Future. *arXiv preprint arXiv:2105.13137* **2021**.
143. Vashishth, S.; Yadati, N.; Talukdar, P. Graph-based deep learning in natural language processing. In *Proceedings of the 7th ACM IKDD CoDS and 25th COMAD*; ACM, 2020; pp. 371–372.
144. Zhang, Z.; Cui, P.; Zhu, W. Deep learning on graphs: A survey. *IEEE Transactions on Knowledge and Data Engineering* **2020**.
145. Lafferty, J.; McCallum, A.; Pereira, F.C. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the 18th International Conference on Machine Learning*; ACM, 2001; pp. 282–289.
146. Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; Bengio, Y. Generative adversarial nets. *Advances in neural information processing systems* **2014**, 27.
147. Kingma, D.P.; Welling, M. An introduction to variational autoencoders. *arXiv preprint arXiv:1906.02691* **2019**.
148. Zhu, D.; Cheng, X.; Zhang, F.; Yao, X.; Gao, Y.; Liu, Y. Spatial interpolation using conditional generative adversarial neural networks. *International Journal of Geographical Information Science* **2020**, 34, 735–758.
149. Li, L.; Franklin, M.; Girguis, M.; Lurmann, F.; Wu, J.; Pavlovic, N.; Breton, C.; Gilliland, F.; Habre, R. Spatiotemporal imputation of MAIAC AOD using deep learning with downscaling. *Remote sensing of environment* **2020**, 237, 111584.
150. Mikołajczyk, A.; Grochowski, M. Data augmentation for improving deep learning in image classification problem. 2018 international interdisciplinary PhD workshop (IIPhDW). IEEE, 2018, pp. 117–122.
151. Zhu, X.X.; Tuia, D.; Mou, L.; Xia, G.S.; Zhang, L.; Xu, F.; Fraundorfer, F. Deep learning in remote sensing: A comprehensive review and list of resources. *IEEE Geoscience and Remote Sensing Magazine* **2017**, 5, 8–36.
152. Castelluccio, M.; Poggi, G.; Sansone, C.; Verdoliva, L. Land use classification in remote sensing images by convolutional neural networks. *arXiv preprint arXiv:1508.00092* **2015**.
153. Penatti, O.A.; Nogueira, K.; Dos Santos, J.A. Do deep features generalize from everyday objects to remote sensing and aerial scenes domains? Proceedings of the IEEE conference on computer vision and pattern recognition workshops, 2015, pp. 44–51.
154. Zhang, S.; Zhang, X.; Zhang, A.; Fu, H.; Cheng, J.; Huang, H.; Sun, G.; Zhang, L.; Yao, Y. Fusion Of Low-And High-Level Features For Uav Hyperspectral Image Classification. 2019 10th Workshop on Hyperspectral Imaging and Signal Processing: Evolution in Remote Sensing (WHISPERS). IEEE, 2019, pp. 1–4.
155. Zou, Q.; Ni, L.; Zhang, T.; Wang, Q. Deep learning based feature selection for remote sensing scene classification. *IEEE Geoscience and Remote Sensing Letters* **2015**, 12, 2321–2325.
156. Chawla, S.; Shekhar, S.; Wu, W.L.; Ozesmi, U. *Modeling spatial dependencies for mining geospatial data: An introduction*; Citeseer, 2000.
157. Schaffer, C. Selecting a classification method by cross-validation. *Machine Learning* **1993**, 13, 135–143.
158. Zheng, M.; Tang, W.; Zhao, X. Hyperparameter optimization of neural network-driven spatial models accelerated using cyber-enabled high-performance computing. *International Journal of Geographical Information Science* **2019**, 33, 314–345.
159. Heremans, S.; Van Orshoven, J. Machine learning methods for sub-pixel land-cover classification in the spatially heterogeneous region of Flanders (Belgium): a multi-criteria comparison. *International Journal of Remote Sensing* **2015**, 36, 2934–2962.
160. Du, Z.; Wu, S.; Zhang, F.; Liu, R.; Zhou, Y. Extending geographically and temporally weighted regression to account for both spatiotemporal heterogeneity and seasonal variations in coastal seas. *Ecological Informatics* **2018**, 43, 185–199.
161. Shi, X.; Yeung, D.Y. Machine learning for spatiotemporal sequence forecasting: A survey. *arXiv preprint arXiv:1808.06865* **2018**.
162. Mazzia, V.; Khaliq, A.; Chiaberge, M. Improvement in land cover and crop classification based on temporal features learning from Sentinel-2 data using recurrent-convolutional neural network (R-CNN). *Applied Sciences* **2020**, 10, 238.
163. Zhang, W.; Thill, J.C. Mesoscale structures in world city networks. *Annals of the American Association of Geographers* **2019**, 109, 887–908.
164. Lee, C.; Wilkinson, D.J. A review of stochastic block models and extensions for graph clustering. *Applied Network Science* **2019**, 4, 1–50.