

# Topo2vec: Topography Embedding Using the Fractal Effect

Jonathan Kavitzky,<sup>1,2</sup> Jonathan Zarecki,<sup>1,3</sup> Idan Brusilovsky,<sup>1,4</sup> Uriel Singer<sup>1,5</sup>

<sup>1</sup> Penta-AI

<sup>2</sup> HUJI: The Hebrew University, Jerusalem, Israel

<sup>3</sup> Bar-Ilan University, Ramat Gan, Israel

<sup>4</sup> The Open University, Ra'anana, Israel

<sup>5</sup> Technion—Israel Institute of Technology, Haifa, Israel

jonathan.kavitzky@mail.huji.ac.il, jonathan.zarecki@gmail.com, brusli1@gmail.com, urielsinger@cs.technion.ac.il

## Abstract

Recent advances in deep learning have transformed many fields by introducing generic embedding spaces, capable of achieving great predictive performance with minimal labeling effort. The geology field has not yet met such success. In this work, we introduce an extension for self-supervised learning techniques tailored for exploiting the fractal-effect in remote-sensing images. The fractal-effect assumes that the same structures (for example rivers, peaks and saddles) will appear in all scales. We demonstrate our method's effectiveness on elevation data, we also use the effect in inference. We perform an extensive analysis on several classification tasks and emphasize its effectiveness in detecting the same class on different scales. To the best of our knowledge, it is the first attempt to build a generic representation for topographic images.

## 1 Introduction

From ancient times, the topographic structure of land was a key aspect in many decisions. The topography of the land is provably correlated with many other tasks: land-use (Sheikh, Van Loon, and Stroosnijder 2014), soil mapping (Scull et al. 2003), soil salinity (Divan and Adriaan 2017), landslides (Prakash, Manconi, and Loew 2020) water floods (Hosseiny, Ghasemian, and Amini 2019), avalanches (Jaedicke, Syre, and Sverdrup-Thygeson 2014) and high solar-energy locations (Heo et al. 2020). The techniques for perceiving, collecting, and understanding topography has changed significantly in recent years and today, geographic information systems (GIS) are built on many classical and data-driven algorithms.

As in many fields in recent years, deep learning has also begun to transform the field of topography and GIS in general, with works from automatic road-mapping (Máttyus, Luo, and Urtasun 2017; Zhang, Liu, and Wang 2018) to elevation map super-resolution (Yue et al. 2015). Deep learning techniques have a particularly interesting property with respect to classical ML models: In addition to superior predictive accuracy, they naturally build a latent embedding space in which every example can be mapped to a latent vector. These embedding spaces have received much attention in self-supervised contrastive learning methods (He et al.

2020; Chen et al. 2020), where an intrinsic property in the data is used to train the model instead of labels gathered by human annotators. Self-supervised methods have been used to achieve close-to-SOTA results on benchmarks such as ImageNet (Russakovsky et al. 2015), iNaturalist (Horn et al. 2018), Birdsnap (Berg et al. 2014), and many more (Nilsback and Zisserman 2008; Fei-Fei, Fergus, and Perona 2004). These methods usually exploit specific invariants in the data, such as invariance to rotation, color jitter, cropping, and more.

Topography data exhibits an interesting property which we call the *fractal-effect*. This effect implies that the same location will have the same topographic patterns when viewed in different radii/scales. Topographic objects such as peaks and saddles will appear in every observed scale. This property also implies in practice that any embedding space for topographic data should be highly scale-dependent, and choosing a scale is highly relevant to how the data would be represented; that is, the same location with different scales should have very different embeddings since they represent different things.

In this work, we present a new self-supervised technique tailored for topographic images, we exploit the fact that topography images (Pelletier 1997) are built as fractals and train a neural network capable of embedding the images in a useful latent embedding space. We use the fractal-effect in inference and prove our embedding space is effective in many scales. Our main contributions are as follows:

1. We introduce a self-supervised training technique for deep neural networks tailored for topography data and use this technique to train a model capable of embedding a topographic image into a useful embedding space.
2. We empirically evaluate the model and technique on classification benchmarks we have built for topographic images, and present qualitative results expressing the fractal-effect in inference.
3. We build a baseline comparison for future work in the field of topography ML. This includes: data, architecture, tasks, and metrics.

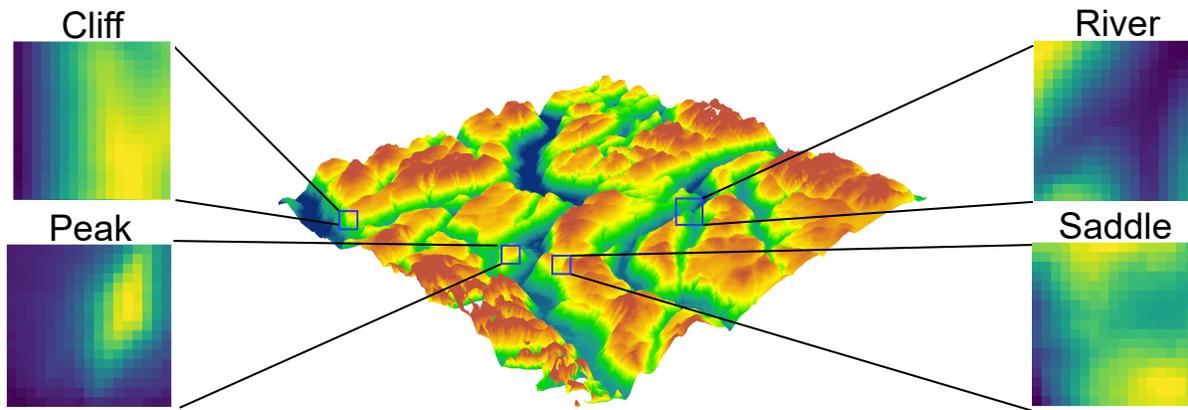


Figure 1: DTM Topography image <sup>1</sup>. Each pixel's value represents its elevation above sea level, and all the values are normalized for each tile.

## 2 Exploiting the Fractal-Effect for Self-Supervision

In this section, we will first introduce the fractal-effect in topography images, afterward, we will discuss the data used in this paper, and then we will define our self-supervision technique exploiting the fractal-effect.

### 2.1 Fractal Effect

As discussed in the introduction, topography data exerts many interesting properties. One of the main properties of the topographic data is scale semi-invariance/semi-fractal: the data can be viewed as topographic in many resolutions and is somewhat similar in those different resolutions (Jin et al. 2017), we call this property the *fractal-effect*. For example, as seen in Figure 2, peaks appear in both very high and low zoom levels of the topography.

As a result of this effect, and in contrast to most computer-vision tasks, our interest objects will certainly appear in multiple scales and any model working in this domain must exert a high level of scale-dependence. Now, in order to support multiple scales, one can change the image resolution so that it will be observed at a different scale. See Figure 2 for a visual illustration.

However, it is important to remember that the fractal-effect implies that each image in the data contains objects in multiple scales. As before performing any experiments one should decide which scale (or multiple scales) the data was sampled at.

### 2.2 Data

**Digital terrain model (DTM)** In order to create the topography images, we used the DTM elevation data from (Agency 2014). Every topography image is a 2D array of pixels, each representing the ground surface elevation above sea level. This database has a spatial resolution of 30 meters per pixel, and elevation accuracy of  $\pm 10$  meters. In our experiments we used Europe as an area of interest.

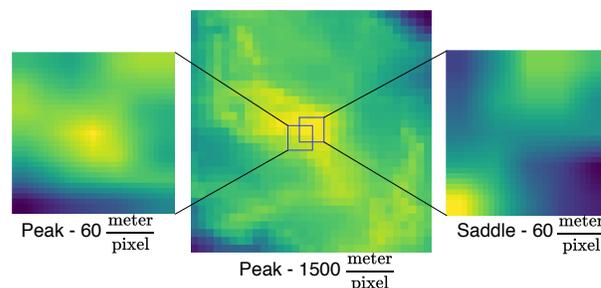


Figure 2: 2 topographic images that were taken from two close points. In a small scale, the two images look very differently and are classified as a peak and a saddle, but in a bigger scale, both locations are of a peak. Different scales imply different topography. Taken in:  $46^{\circ}51'34.2'' N$   $7^{\circ}31'31.8'' E$ .

**OpenStreetMap (OSM)** For precise locations of objects of different classes, used later for evaluation, we used OpenStreetMap (Haklay and Weber 2008). Each location is represented by a pair of geographic coordinates (longitude, latitude), which in turn will be converted to a topography image. We used geo-locations of different topographic classes (peaks, rivers and saddles), and non-topographic classes (aerialway stations) all sampled evenly at random out of all the OSM known tags.

### 2.3 Topo2Vec

Let  $loc$  be a location on earth:  $x_r^s$  is the image of the area with radius  $r$  around  $loc$  (as presented in Figure 1),  $N$  is the number of pixels from the image's central to the edge. The spatial resolution in  $\frac{\text{meter}}{\text{pixel}}$  of  $x_r^s$  is defined as:

$$s = r/N \frac{\text{meter}}{\text{pixel}}$$

One can sample the image at different spatial resolutions  $x_r^{s_1}, x_r^{s_2}, \dots, x_r^{s_j}$ , while covering the same area. However, due to the fractal-effect, different topographic objects in the

<sup>1</sup>gisgeography.com

observed area will become more apparent in different resolutions, as seen in Figure 2.

Taking the fractal-effect into consideration, we define  $f$  as a function that encodes  $x_r^s$  into an embedding space. Our goal is that  $f(x_r^s)$  will represent the topographic information present around  $loc$  with radius  $r$  and a resolution of  $s$ . By learning a second function,  $g$ , that decodes  $f(x_r^s)$  to an image with a different spatial resolution  $s'$ , we are able to leverage the fractal-effect and learn a better topographic representation. This is achieved by maximizing:

$$\mathbb{P}(x_r^{s'} | g(f(x_r^s)))$$

Where  $s' = k \cdot s$  ( $k \geq 1$ ) can be thought as the *fractal-factor*.

Maximizing the above expression is done by minimizing the  $L^p$  distance:

$$\min_{\theta_f} L_p = \|x_r^{s'}, g(f(x_r^s))\|$$

Training using this loss, we obtain two functions  $f$  and  $g$ . While  $g$  decodes a given embedding vector into a higher resolution image, meaning it learned the fractal-effect around  $loc$ ,  $f$  (from now on will be called *topo2vec*) is a self-supervised, fractal-effect aware encoder that takes as input a topographic image and embeds it into an embedding space.

---

**Algorithm 1:** Topo2vec training procedure

---

**Input :** *Locations* - list of different locations

*Scales* - list of relevant scales

$k$  - desired fractal-factor

$lr_1, lr_2$  - learning rates

**Output:**  $f$  - learned encoder

$g$  - learned decoder

$d$  - learned discriminator

```

1 while  $f, g, d$  not converged do
2   for  $batch$  in Locations do
3     for  $s \in$  Scales do
4        $X = \text{toImageDataset}(batch, s)$ ;
5        $Y = \text{toImageDataset}(batch, k \cdot s)$ ;
6        $\hat{Y} = g(f(X))$ ;
7        $L_p = \sum_i \|y_i, \hat{y}_i\|_p$ ;
8        $L_G^{cGAN} = BCE(d([X, \hat{Y}], \bar{1}))$ ;
9        $f \leftarrow f - lr_1 \cdot \nabla_f(\lambda_1 * L_p + \lambda_2 * L_G^{cGAN})$ ;
10       $g \leftarrow g - lr_1 \cdot \nabla_g(\lambda_1 * L_p + \lambda_2 * L_G^{cGAN})$ ;
11       $L_D^{cGAN} = BCE(d([X, Y], \bar{1})) +$ 
            $BCE(d([X, \hat{Y}], \bar{0}))$ ;
12       $d \leftarrow d - lr_2 \cdot \nabla_d(L_D^{cGAN})$ ;
13 return  $f, g, d$ ;
```

---

Intuitively, if  $loc$  expresses a fractal-effect that one peak becomes many small peaks in higher resolution. Then, if  $g$  has managed to predict this fractal-pattern, it has learned the topographic information in  $loc$ . A problem that arises from this training procedure is that while  $g$  attempts to predict the fractal-pattern, it is not guaranteed that it will predict the exact fractal-effect of  $loc$ , but rather the distribution of patterns possible in that area. In order to address this issue, we

propose an adversarial loss on top of the current loss to enable the network to predict a sample from the fractal-pattern distribution rather than the distribution's average. The discriminator  $d$  takes as input fake pairs,  $(x_r^s, g(f(x_r^s)))$ , and true pairs,  $(x_r^s, x_r^{s'})$  and tries to distinguish between them, while  $g \circ f$  tries to fool him. Formally, the adversarial loss is:

$$\min_{\theta_f \theta_g} \max_{\theta_d} L_G^{cGAN} = \mathbb{E}_{s \sim (x_r^s, x_r^{s'})} \log(d(s)) \\ + \mathbb{E}_{s \sim (x_r^s, g(f(x_r^s)))} \log(1 - d(s))$$

While the final optimization loss of topo2vec is:

$$\min_{\theta_f \theta_g} \max_{\theta_d} \lambda_1 * L_p + \lambda_2 * L_G^{cGAN}$$

Where  $\lambda_1 = 100$  and  $\lambda_2 = 1$ , and  $L_G^{cGAN}$  is calculated using binary cross-entropy (*BCE*). Alg 1 presents the training procedure of the adversarial version of topo2vec, while non-adversarial topo2vec trains the same but with  $\lambda_1 = 1$  and  $\lambda_2 = 0$ .

## 2.4 Implementation Details

The topo2vec architecture consists of down-sample layers and then up-sample layers, similar to the U-net architecture (Ronneberger, Fischer, and Brox 2015), but without the skip connections. We take the representation between the down and up-sample layers as the latent representation. The reason we deleted the skip connections is to constrain the latent representation to hold the entire information needed for decoding. In order to allow learning  $s' > s$  we added additional up-sample layers (one for each factor of 2). We used L1 distance as the loss with a learning rate of  $lr_1 = 0.0002$  (for the adversarial version  $g$  had the learning rate of  $lr_2 = 0.0016$ ). As for the discriminator, we trained a 5-layer CNN with three fully-connected layers at the end.

**Encoder** The full architecture of the encoder  $f$  is as follows:

1. The input is a 1x17x17 image representing the topography of a specific coordinate in a given scale.
2. The input image enters two convolutional layers, each includes a convolution block (8 kernels of size 3x3), BatchNorm (Ioffe and Szegedy 2015), and ReLu (Nair and Hinton 2010). The image is resized to 16x16 (due to padding), with 8 channels, resulting in an 8x16x16 output.
3. Four down-sample layers, each includes a max-pooling layer (kernel size of 2x2) followed by two convolutional layers as before, each time doubling the number of filters. After four down-sample layers, we end up with a 128x1x1 image.

The 128x1x1 image is flattened into a 128 vector, representing the latent vector of the image.

**Decoder** The full architecture of the decoder  $g$  is as follows:

1. The input is a 128x1x1 image representing the last output of  $f$ .

- Four up-sample layers, each includes an up-sampling layer (scale factor of 2 in bilinear mode), followed by two convolutional layers as before, each time decreasing the number of filters by 2. After 4 up-sample layers, we end up with an  $8 \times 16 \times 16$  image.
- In `topo2vec-1` we finish with a convolutional block that maps the  $8 \times 16 \times 16$  image to the final  $1 \times 16 \times 16$  output image.
- In `topo2vec-4` we add two additional up-sample layers that resolve with a  $2 \times 64 \times 64$  image. We finish with a convolutional block that maps the  $2 \times 64 \times 64$  image to the final  $1 \times 64 \times 64$  output image.

**Discriminator** For the `topo2vec-adv` version,  $f$  and  $g$  together represent the generator, while  $d$  represent the discriminator. The full architecture of the discriminator  $d$  is as follows:

- The input is a  $2 \times 64 \times 64$  image representing in its two channels the true or fake pair of images.
- Five layers of convolutional blocks (For the first four: kernel size of 4, stride 2, and padding 1. For the Fifth: kernel size of 4, stride 1, and no padding), BatchNorm, and ReLu. The first starts with 8 filters and each time doubling, ending up with a  $128 \times 1 \times 1$  image.
- Flattening the image to a 128 size vector.
- Three fully connected layers:  $128 \rightarrow 64, 64 \rightarrow 32, 32 \rightarrow 1$ , with a ReLu activation after the first two.
- Sigmoid activation resolving with a final probability neuron.

### 3 Experiments

We test our framework on several topographic classes gathered from OpenStreetMap (OSM), the leading geodata repository, in both few-shot and many-shot settings. Lastly, we'll show an experiment to show how we generalize beyond our initially trained scale.

#### 3.1 Experimental Methodology

**Datasets** In order to guarantee a fair evaluation procedure, we geographically split the DTM into two different areas of size  $25 \text{ degree}^2$  each, one for train and one for test (Figure 3). The polygon from which we sampled coordinates for the train-set is:

$POLYGON((10 \ 50, 10 \ 45, 15 \ 45, 15 \ 50, 10 \ 50))$

The polygon from which we sampled coordinates for the test set is:

$POLYGON((5 \ 45, 5 \ 50, 10 \ 50, 10 \ 45, 5 \ 45))$

Each sample was built by extracting topographic data on an image of radius  $r$  around a geographic point (as presented in Figure 1), and resizing it to the appropriate size ( $17 \text{ pixels}$ ).

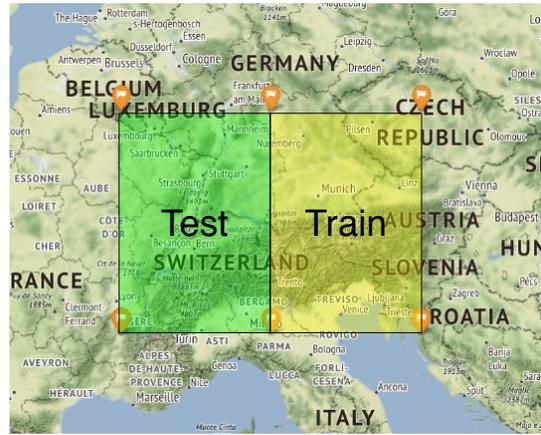


Figure 3: Geographic train/test split.

#### Baselines

- ID**: A simple flatten operation on the given image resulting in a vector that is used as the feature representation.
- CNN**: A convolutional network that is trained on the classification of geographic entities from OSM (fully-supervised learning). The CNN is made of 2 convolutional layers and 2 fully-connected layers. The model is trained on examples of peaks, rivers, cliffs, and saddles - the same classes we later examine in Exp. 1, as those are the only "basic" topographic classes with an appropriate amount of data in the OSM for deep networks training, this results in 20,000 training examples. We use the last hidden layer as the feature representation.
- SauMoCo** (Kang et al. 2020): An extension of the contrastive method MoCo (He et al. 2020). This method uses an additional spatial augmentation technique of sampling a geographically close image. We use this method instead of classic MoCo as other augmentations employed by contrastive methods are irrelevant for topography data, this includes color jitter, random conversion to greyscale, and more. This method replaces using Tile2Vec (Jean et al. 2019) as it clearly outperforms it in spatial image representation.

We additionally compare to several variations of our method:

- Topo2vec-1**: Topo2vec with fractal-factor  $k = 1$  (meaning no fractal learning). This variation will help us compare the effectiveness of exploiting the fractal-effect.
- Topo2vec-4**: Topo2vec with fractal-factor  $k = 4$ .
- Topo2vec-adv**: An adversarial version of topo2vec network as explained last section with fractal-factor  $k = 4$ .

All self-supervised methods were trained with 100,000 images. All input images to all methods are the size of  $17 \times 17^2$ .

<sup>2</sup>GitHub repository with all code, baselines, data, and experiments: <https://github.com/urielsinger/topo2vec>

### 3.2 Exp. 1: Classification on Topographic Classes

In this experiment, we evaluate a model trained by our framework’s ability to serve as a pre-training model for topographic tasks with a sufficient amount of data. We test our model individually on 4 topographic classes (rivers<sup>3</sup>, peaks<sup>4</sup>, saddles<sup>5</sup> and cliffs<sup>6</sup>).

**Finding the class’ OSM scale** A given topographic class is present at many scales, but in OSM, not all scales were collected. We’ll first devise an experiment to determine what scale the OSM community sampled at:

Let  $C = \{c \mid c \in \text{coordinates}\}$  be a coordinate dataset for a given class, constructed such that half of the coordinates are sampled at the location of the class ( $y = 1$ ) and half are sampled at random ( $y = 0$ ). Given a radius  $r$ , we define  $X_r^C = \{x_r^i \mid x_r^i \in \mathbb{R}^{17 \times 17}\}$  as the image dataset built from  $C$  where each image is built around a coordinate with radius  $r$ . For each  $r$  the spatial resolution changes.

Now, to find the sampled resolution/scale  $s$  of set  $C$  we define the following experiment:

1. Split  $X_r^C$  of size 1,000 to train and validation sets (80%, 20%).
2. For each  $r \in R$ , train a CNN on the available training data and calculate accuracy on the validation set. This CNN has a similar architecture to the CNN baseline.
3. Repeat this process for 10 different random seeds
4. The sampled scale is the one with maximum average accuracy on the validation set.

We can expect the model to have maximum accuracy on the scale of the class’ sampled scale, the experiment is repeated for all four classes. The results of the experiment classes are presented in Figure 4. We observe that all classes peak around the resolution of 30 meters/pixel. All subsequent experiments will train the models in this optimal scale.

**Methodology & Results** We compared the topo2vec variants with baselines trained on the optimal scale presented above. These models are pre-trained as discussed in 3.1.

For each topographic class and model, we train an SVM built on the model’s embeddings with a training set of size 1000 and evaluated on a test set of size 200, both equally distributed between the positive and negative examples. We repeat each experiment 10 times with random seeds and report average accuracies and standard deviations in Table 1.

We see that our model outperforms other baselines on all classes, even when the baselines are trained on the optimal scale. It is important to notice that although CNN is a supervised baseline which had access to more labeled data, topo2vec outperformed it over all classes. Furthermore, topo2vec-4 outperforms topo2vec-1 on 2 classes and is comparable on the other two. This shows the power the fractal-effect has in the self-supervised learning procedure. Topo2vec-adv outperforms topo2vec-4 on two classes and

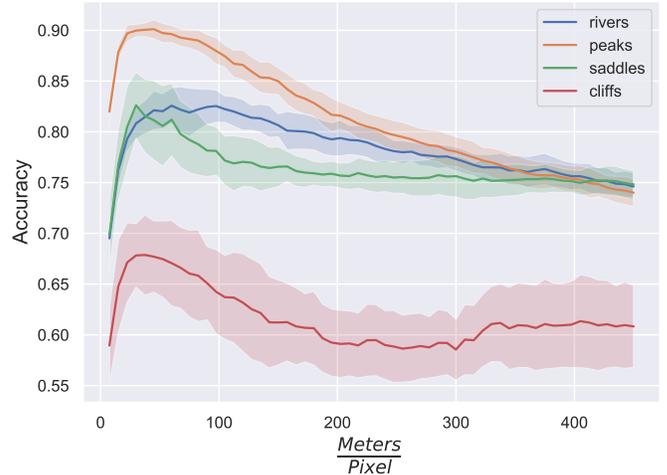


Figure 4: Accuracy-per-scale of the topography classes.

	peaks	rivers	cliffs	saddles
ID	92.4 ± 0.5	75.9 ± 1.8	74.4 ± 2.8	87.5 ± 1.2
CNN	92.3 ± 0.5	76.4 ± 0.7	72.2 ± 1.8	84.1 ± 1.3
SauMoCo	86.9 ± 2.3	81.7 ± 2.2	66.2 ± 2.5	67.3 ± 2.3
topo2vec-1	<b>93.5 ± 0.7</b>	81.3 ± 2.2	71.4 ± 1.7	80.5 ± 2.3
topo2vec-4	93.1 ± 0.9	80.9 ± 1.7	75.8 ± 2.1	<b>89.0 ± 1.3</b>
topo2vec-adv	92.8 ± 0.3	<b>82.6 ± 1.2</b>	<b>78.1 ± 2.8</b>	88.7 ± 2.2

Table 1: Experiment 1 results. accuracy ± standard-deviation

shows the same results on the other two. This indicates that the adversarial loss affects the latent representation learning, although not significantly. SauMoCo does not perform well in this experiment, we hypothesize this is because several of its base assumptions does not hold in topography data. It assumes that close locations should have the same representation (similar to Tile2Vec (Jean et al. 2019)), and in topography data, where the fractal-effect is present, even the same location can have multiple classes when viewed in different scales. While SauMoCo did not deal with this fact, it is the key consideration of our method.

### 3.3 Exp. 2: Classification on Topography-correlated Classes

We evaluate our method on several topography-correlated classes, where we want to see if the embedding space holds information useful even for classes with only a mild topographic signature.

We test on four classes: aerialway station<sup>7</sup>, alpine hut<sup>8</sup>, waterfall<sup>9</sup> and sinkhole<sup>10</sup>. Similar to Exp. 1, we train an SVM built on the model’s embeddings. We sample a training set of size 400 and a test set of size 100, both equally

<sup>3</sup>wiki.openstreetmap.org/wiki/Rivers

<sup>4</sup>wiki.openstreetmap.org/wiki/Peaks

<sup>5</sup>wiki.openstreetmap.org/wiki/Tag:natural=saddle

<sup>6</sup>wiki.openstreetmap.org/wiki/Tag:natural=cliff

<sup>7</sup>wiki.openstreetmap.org/wiki/Key:aerialway

<sup>8</sup>wiki.openstreetmap.org/wiki/Tag:tourism=alpine\_hut

<sup>9</sup>wiki.openstreetmap.org/wiki/Waterfalls

<sup>10</sup>wiki.openstreetmap.org/wiki/Tag:natural=sinkhole

distributed between the positive and negative examples. We repeat each experiment 10 times with random seeds. The average accuracies and standard deviations are summarised in Table 2.

	aerialway stations	alpine huts	waterfall	sinkholes
ID	63.9 ± 5.5	59.2 ± 3.6	69.2 ± 2.6	62.5 ± 2.4
CNN	75.5 ± 2.1	73.3 ± 1.6	69.7 ± 2.4	55.3 ± 2.9
SauMoCo	69.7 ± 3.6	73.3 ± 2.8	55.3 ± 3.1	56.5 ± 4.3
topo2vec-1	79.4 ± 1.5	70.9 ± 2.7	<b>77.9 ± 2.2</b>	60.1 ± 2.1
topo2vec-4	82.1 ± 2.0	<b>75.7 ± 1.4</b>	73.3 ± 1.6	<b>63.2 ± 2.5</b>
topo2vec-adv	<b>82.8 ± 1.8</b>	75.5 ± 1.3	74.8 ± 1.6	62.9 ± 2.4

Table 2: Experiment 2 results. accuracy ± standard-deviation

We can see how topo2vec significantly outperforms the other baselines on this task. This result emphasizes, wherein the previous experiment all methods were able to learn the topographic classes, it is much harder to predict classes with a less emphasized topographic signature. topo2vec-4 outperforms topo2vec-1 over all classes except waterfall, topo2vec-adv show comparable results to topo2vec-4. The poor performance in the waterfall class can be explained as waterfall do not exhibit fractal properties (no peak inside a waterfall).

This experiment shows that our modifications for topographic data are useful for learning a more powerful representation, as intended. The exploitation of the fractal-effect improves the generality of topography representation that was learned. Moreover, this suggests that topography data holds additional, uncorrelated information that can be used in addition to normal GIS and tabular features to further improve ML models working in the GIS domain.

### 3.4 Exp. 3: Fractal-Effect Visualization

Throughout this work, we repeatedly presented the importance of the fractal-effect and how we might find entirely different objects when looking at different scales. In this section, we will present several qualitative results of topo2vec-4 exhibiting this important property.

For each point in the polygon, we build an image around it, which is then classified using our model. This process is repeated for different radii and scales. In Figure 5 we can see our method detecting peaks, saddles, rivers, and cliffs in multiple scales when images in different zoom levels are presented. This is a unique property of topography data and we see our model performs well in this setting.

### 3.5 Exp. 4: Topography Retrieval

We conducted another qualitative experiment to demonstrate the strength of our embedding space. We took  $n$  locations representing a topography pattern as input and averaged over their topo2vec-4 embedding representations. We then used the KNN algorithm to find the “most similar” points in the embedding space to those input points. In Figure 6 we can see a representative example. We can see that our retrieved

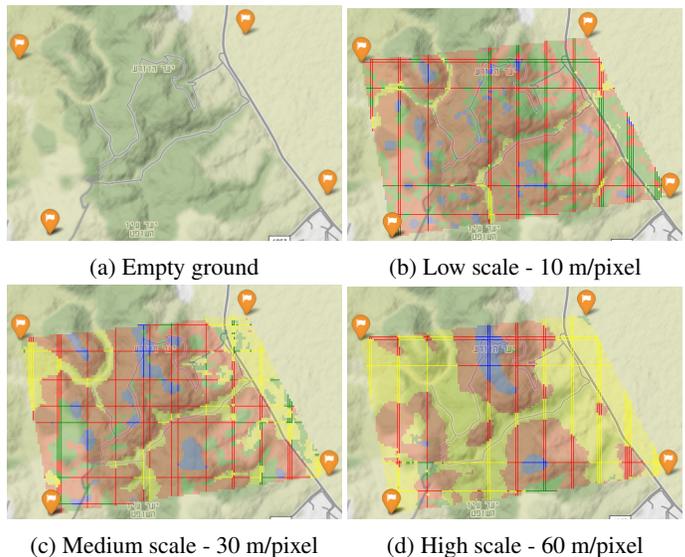


Figure 5: Qualitative results of our model’s predictions in different scales. (blue: peaks, green: saddles, red: cliffs, yellow: rivers) This polygon is around  $32^{\circ}37'02.8'' N 35^{\circ}07'27.4'' E$ .

coordinates have very similar topography-patterns to those sampled.

This interactive experiment is also provided in the code.

## 4 Related Work

### 4.1 Algorithmic Topography

There have been works dealing with topography data enhancement and generation: Extracting information from raw elevation data (Jaedicke, Syre, and Sverdrup-Thygeson 2014), DTM (digital-terrain-model) extraction from satellite imagery (Gevaert et al. 2018) and enhancing the resolution of topographic images (Yue et al. 2015). Many possible usages of DTM have been proposed. natural disasters analysis and prediction (such as Avalanche warning (Jaedicke, Syre, and Sverdrup-Thygeson 2014; Choubin et al. 2019) and landslide susceptibility (Wu and Chen 2009)) to high solar energy regions localization (Heo et al. 2020) and more (Bolibar et al. 2020). In addition, deep learning has been used for the automatic mapping of topographic objects from DTM (Torres et al. 2018; Torres, Milani, and Fraternali 2019), and satellite imagery (Li and Hsu 2020).

### 4.2 Fractal-Effect in Topography

The fact the earth’s topography is a fractal has inspired some previous works. (Chase 1992) studied the evolution of mountains and regional topography, and the effects of tectonic movement and climate on the landscape, including its fractal geometry. (Pelletier 1997) built a mathematical framework for explaining the existence of fractals in topography. (Weissel et al. 1995) used the fractal-effect in a research of the erosional development of the Ethiopian plateau

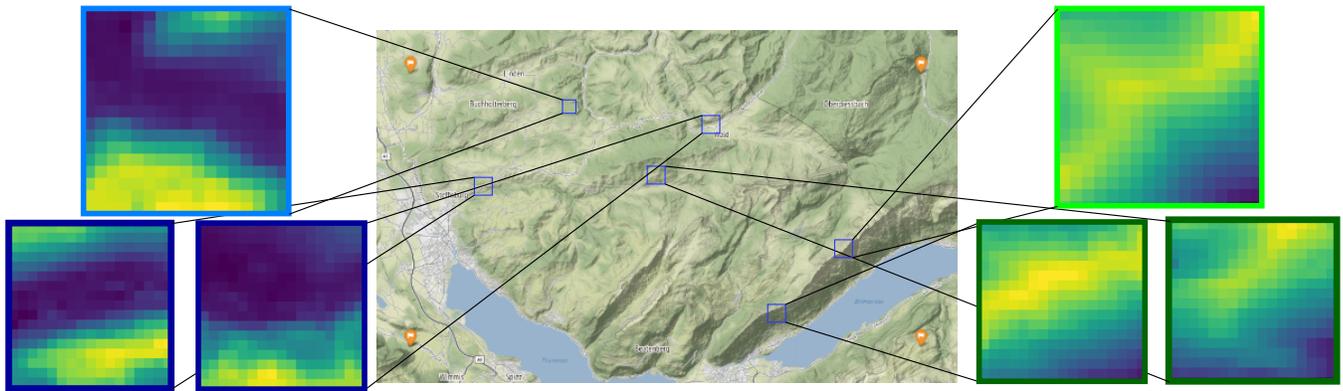


Figure 6: The top image of each side is the input location. The two bottom images are the two nearest neighbors of that input image in the latent space. This experiment was conducted in the orange rectangle, around  $46^{\circ}46'30.0'' N 7^{\circ}51'00.0'' E$ .

of Northeast Africa and (Liu et al. 2019) used it in a similar way for landslide susceptibility mapping.

### 4.3 Unsupervised Learning for Visual and Geographic Data

The concept of embedding an input to a latent lower-dimensional space is a basic and powerful one in the field of machine learning. As mentioned in the introduction, neural networks tend to generate these spaces naturally during training. Unsupervised learning for visual tasks is a very active field of research, making it very difficult to summarize. Recent methods that have pushed the state-of-the-art include: SimCLR (Chen et al. 2020), MoCo (He et al. 2020) and BYoL (Grill et al. 2020). All these methods follow the basic idea of making representation of an image similar under small transformations, exploiting several invariances within their domain such as invariance to rotation and color jitter. There have been some attempts in the geographic ML field to build specialized embedding spaces, most following Tobler’s first law of geography (Tobler 1970): “*everything is related to everything else, but near things are more related than distant things*”. Noteworthy examples are Tile2Vec (Jean et al. 2019) which used the triplet-loss for spatially distributed data and the recent work of SauMoCo (Kang et al. 2020) which takes this same basic idea from Tobler’s law in order to create a new invariance for *spatially close images*. We have seen in our experiments this assumption does not hold for classification on topographic classes.

## 5 Conclusions

In this work, we presented a novel self-supervision method for training neural networks for topographic feature extraction. We exploited the fractal-effect in the data during training and inference to build a model capable of generalizing to any relevant scale. Our key consideration in topo2vec was leveraging the fractal-effect, which we achieved in an encoder-decoder based framework. We evaluated our method on several topographic classes and topography-correlated classes and found it superior to other

self-supervised methods and even surpassing fully supervised methods that used an order of magnitude more data.

Our results motivate a few interesting directions for future work we intend to explore. First is pushing empirical results further by using more elaborate architectures such as graph autoencoders (Kipf and Welling 2016) and GANs (Goodfellow et al. 2014). Second, we plan to use our method for the automation of topography features mapping in multiple scales, greatly reducing the human labor necessary for such a task.

We hope this work will further advance the field of geographic ML, specifically to be used as additional features in tasks that exhibit correlation with the topography of its location (like avalanches and wildfire prediction). We also hope our method will act as a strong baseline for future works in the field of topography embedding.

## Bibliography

- Agency, J. A. E. 2014. Advanced Land Observing Satellite DAICHI-2. DAICHI.
- Berg, T.; Liu, J.; Lee, S.; Alexander, M. L.; Jacobs, D.; and Belhumeur, P. 2014. Birdsnap: Large-Scale Fine-Grained Visual Categorization of Birds. *CVPR*, 2019–2026.
- Bolibar, J.; Rabatel, A.; Gouttevin, I.; Galiez, C.; Condom, T.; and Eric, S. 2020. Deep learning applied to glacier evolution modelling. *The Cryosphere*, 14: 565–584.
- Chase, C. G. 1992. Fluvial land sculpting and the fractal dimension of topography. *Geomorphology*, 5(1): 39 – 57. *Fractals in Geomorphology*.
- Chen, T.; Kornblith, S.; Norouzi, M.; and Hinton, G. E. 2020. A Simple Framework for Contrastive Learning of Visual Representations. *ArXiv*, abs/2002.05709.
- Choubin, B.; Borji, M.; Mosavi, A.; Sajedi-Hosseini, F.; Singh, V. P.; and Shamshirband, S. 2019. Snow avalanche hazard prediction using machine learning methods. *Journal of Hydrology*, 577: 123929.
- Divan, V.; and Adriaan, V. N. 2017. Machine learning performance for predicting soil salinity using different combinations of geomorphometric covariates. *Geoderma*, 299: 1 – 12.

- Fei-Fei, L.; Fergus, R.; and Perona, P. 2004. Learning Generative Visual Models from Few Training Examples: An Incremental Bayesian Approach Tested on 101 Object Categories. In *CVPR Workshops*.
- Gevaert, C.; Persello, C.; Nex, F.; and Vosselman, G. 2018. A deep learning approach to DTM extraction from imagery using rule-based training labels. *ISPRS Journal of Photogrammetry and Remote Sensing*, 142.
- Goodfellow, I. J.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A. C.; and Bengio, Y. 2014. Generative Adversarial Nets. In *NIPS*.
- Grill, J.-B.; Strub, F.; Altché, F.; Tallec, C.; Richemond, P. H.; Buchatskaya, E.; Doersch, C.; Pires, B. A.; Guo, Z. D.; Azar, M. G.; Piot, B.; Kavukcuoglu, K.; Munos, R.; and Valko, M. 2020. Bootstrap Your Own Latent: A New Approach to Self-Supervised Learning. *ArXiv*, abs/2006.07733.
- Haklay, M.; and Weber, P. 2008. Openstreetmap: User-generated street maps. *IEEE Pervasive Computing*, 7(4).
- He, K.; Fan, H.; Wu, Y.; Xie, S.; and Girshick, R. B. 2020. Momentum Contrast for Unsupervised Visual Representation Learning. *CVPR*, 9726–9735.
- Heo, J.; Jung, J.; Kim, B.; and Han, S. 2020. Digital elevation model-based convolutional neural network modeling for searching of high solar energy regions. *Applied Energy*, 262: 114588.
- Horn, G. V.; Aodha, O. M.; Song, Y.; Cui, Y.; Sun, C.; Shepard, A.; Adam, H.; Perona, P.; and Belongie, S. J. 2018. The iNaturalist Species Classification and Detection Dataset. *CVPR*, 8769–8778.
- Hosseiny, B.; Ghasemian, N.; and Amini, J. 2019. A convolutional neural network for flood mapping using sentinel-1 and STRM DEM data: case study in Poldokhtar-Iran. volume XLII-4/W18, 527–533.
- Ioffe, S.; and Szegedy, C. 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*.
- Jaedicke, C.; Syre, E.; and Sverdrup-Thygeson, K. 2014. GIS-aided avalanche warning in Norway. *Computers & Geosciences*, 66: 31 – 39.
- Jean, N.; Wang, S.; Samar, A.; Azzari, G.; Lobell, D.; and Ermon, S. 2019. Tile2Vec: Unsupervised representation learning for spatially distributed data. In *AAAI*.
- Jin, Y.; Wu, Y.; Li, H.; Zhao, M.; and Pan, J. 2017. Definition of fractal topography to essential understanding of scale-invariance. *Scientific Reports*, 7: 46672.
- Kang, J.; Fernandez-Beltran, R.; Duan, P.; Liu, S.; and Plaza, A. J. 2020. Deep Unsupervised Embedding for Remotely Sensed Images Based on Spatially Augmented Momentum Contrast. *IEEE Transactions on Geoscience and Remote Sensing*, 1–13.
- Kipf, T. N.; and Welling, M. 2016. Variational Graph Auto-Encoders. *NIPS Workshop on Bayesian Deep Learning*.
- Li, W.; and Hsu, C. 2020. Automated terrain feature identification from remote sensing imagery: a deep learning approach. *International Journal of Geographical Information Science*, 34: 637 – 660.
- Liu, L.; Li, S.; Li, X.; Jiang, Y.; Wei, W.; Wang, Z.; and Bai, Y. 2019. An integrated approach for landslide susceptibility mapping by considering spatial correlation and fractal distribution of clustered landslide data. *Landslides*, 16: 715–728.
- Máttyus, G.; Luo, W.; and Urtasun, R. 2017. Deep-RoadMapper: Extracting Road Topology from Aerial Images. *ICCV*, 3458–3466.
- Nair, V.; and Hinton, G. E. 2010. Rectified linear units improve restricted boltzmann machines. In *ICML*.
- Nilsback, M.-E.; and Zisserman, A. 2008. Automated Flower Classification over a Large Number of Classes. *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*, 722–729.
- Pelletier, J. D. 1997. Why is topography fractal.
- Prakash, N.; Manconi, A.; and Loew, S. 2020. Mapping landslides on EO data: Performance of deep learning models vs. traditional machine learning models. *Remote Sensing*, 12(3): 346.
- Ronneberger, O.; Fischer, P.; and Brox, T. 2015. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, 234–241. Springer.
- Russakovsky, O.; Deng, J.; Su, H.; Krause, J.; Satheesh, S.; Ma, S.; Huang, Z.; Karpathy, A.; Khosla, A.; Bernstein, M.; Berg, A.; and Fei-Fei, L. 2015. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 115: 211–252.
- Scull, P.; Franklin, J.; Chadwick, O.; and McArthur, D. 2003. Predictive soil mapping: a review. *Progress in Physical Geography*, 27(2): 171–197.
- Sheikh, V.; Van Loon, E.; and Stroosnijder, L. 2014. Relationship between topography, land use and soil moisture in loess hillslopes. *Environmental Resources Research*, 1(2): 141–165.
- Tobler, W. R. 1970. A Computer Movie Simulating Urban Growth in the Detroit Region. *Economic Geography*, 46: 234–240.
- Torres, R. N.; Fraternali, P.; Milani, F.; and Frajberg, D. 2018. A Deep Learning Model for Identifying Mountain Summits in Digital Elevation Model Data. *AIKE*, 212–217.
- Torres, R. N.; Milani, F.; and Fraternali, P. 2019. Algorithms for Mountain Peaks Discovery: A Comparison. In *ACM/SIGAPP, SAC '19*, 667–674. New York, NY, USA: Association for Computing Machinery. ISBN 9781450359337.
- Weissel, J. K.; Malinverno, A.; Harding, D. J.; and Karner, G. D. 1995. *Erosional Development of the Ethiopian Plateau of Northeast Africa from a Fractal Analysis of Topography*, 127–142. Boston, MA: Springer US. ISBN 978-1-4615-1815-0.
- Wu, C.-H.; and Chen, S.-C. 2009. Determining landslide susceptibility in Central Taiwan from rainfall and six site factors using the analytical hierarchy process method. *Geomorphology*, 112(3): 190 – 204.
- Yue, L.; Shen, H.; Yuan, Q.; and Zhang, L. 2015. Fusion of multi-scale DEMs using a regularized super-resolution

method. *International Journal of Geographical Information Science*, 29: 2095 – 2120.

Zhang, Z.; Liu, Q.; and Wang, Y. 2018. Road Extraction by Deep Residual U-Net. *IEEE Geoscience and Remote Sensing Letters*, 15: 749–753.