Masked Image Modelling for Representation Learning in Earth Observation

Hugo Hernández Hernández

Helmholtz AI Conference 2023

Hamburg, 12.06.2023

#### Contents

- Introduction
- Methodology
- Experiments
- Results
- Conclusions

# Introduction

# 1. Introduction | Background | Self-supervised learning

- SSL aims to learn by itself how to learn massive amount of data
- Masked Image Modeling as part of SSL



(Wang et al., 2022)

# Methodology

## 2. Methodology | MIM pre-training | Masked Autoencoders



(He et al., 2021)

#### 2. Methodology | MIM pre-training | Masked Feature Prediction



# 2. Methodology | MIM pre-training | MAE + MFP



# 2. Methodology | MIM pre-training | Histogram Oriented Gradients





Representation of gradients

#### Visualization of HOG

# Experiments

# 3. Experiments | EO datasets

#### SSL4EO-s12

- 3 million Sentinel-2 and sentinel-1 images
- 250k locations sampled
- Patch scale 264x264 px
- Multi-spectral images over 13 bands

#### EuroSAT

- 27k labeled and georeferenced images
- Locations distributed all over Europe
- Patch scale 64x64 px
- Multi-spectral images over 13 bands

# 3. Experiments | Self-supervised pre-training



# 3. Experiments | Downstream Task



# 3. Experiments | Ablation studies

Masking ratio



Normalization

RAW HOG



# 3. Experiments | Evaluation strategies



Supervised Learning

• Random Initialization

• Labels variation

# Results

# 4. Results | SSL pre-training | Image reconstruction



#### 4. Results | Downstream task | Classification accuracies



#### 4. Results | Downstream task | Evaluation strategies



■ Random Initialization ■ Supervised Learning ■ Linear classification(RAW) ■ Linear classification (HOG) ■ Fine Tuning (RAW) ■ Fine Tuning(HOG)

#### 4. Results | Downstream task | Normalization



### 4. Results | Downstream task | Masking ratio



## 4. Results | Downstream task | Confusion Matrix 100%



## 4. Results | Downstream task | Misclassified images





(b)



(c)

(a)



(d)



(e)

## 5. Conclusions and Outlook

- Hybrid model performed in MSI remote sensing using SSL
- Classification accuracy surpasses state of the art, feature descriptors slightly improve performance
- Performance do not present a big improvement with respect pixelwise analysis
- Open possibility to explore new feature descriptors

#### Thanks for your attention! Hugo Hernández Hernández

# Appendix

## Appendix | Motivation

- Human annotated datasets are costly and time-consuming
- Earth observation deals with massive-scale datasets, self-supervised learning can address this issue

# **Appendix | Vision Transformers**

Vision Transformer (ViT)



(Dosovitskiy et al., 2020)

# Appendix | Downstream task



#### Linear Classification

Fine-tuning





(Dosovitskiy et al., 2020)

# Appendix | Data preparation



- 27k images
- 10 categories

# Appendix | Ablation study





(Towards Al, 2020)

Masking ratio

Normalization

# Appendix | SSL pre-training | Training loss

Mask/Norm	Α	В	С	D	Ε	F
0.7	4.55E-04	0.038	0.664	0.034	5.74E-06	5.75E-06
0.5		2.78E-04	0.593	0.030		
0.2		1.72E-04	0.525	0.026		

- A RAW normalized
- B RAW no normalized
- C HOG normalized + hog-normalized

- D HOG hog-normalized
- E-HOG normalized
- F HOG no normalized

### Appendix | Downstream task | Evaluation strategies

		100%	50%	10%	1%
<b>Random Initialization</b>		0.7840	0.7929	0.795	0.7062
Supervised Learning		0.969	0.9459	0.8729	0.6031
Linear	RAW	0.9267	0.9367	0.8935	0.8459
Classif.	HOG	0.9196	0.9398	0.9187	0.8798
Fine	RAW	0.9857	0.9742	0.9405	0.7998
Tuning	HOG	0.9872	0.9796	0.9442	0.8053

#### Appendix | Downstream task | Normalization

		norm	norm.hog	n+norm.hog	
Linear	HOG	0.8598	0.9398	0.9125	0.8794
Classif.					
Fine	HOG	0.9661	0.9824	0.9872	0.9701
Tuning					

### Appendix | Downstream task | Masking ratio

		M: 0.7	M: 0.5	M: 0.2
Linear	HOG	0.9196	0.9277	0.9283
Classif.				
Fine	HOG	0.9872	0.9838	0.9877
Tuning				


# 4. Results | Downstream task | Confusion Matrix 10%



# Appendix | Vision Transformer complete





(Dosovitskiy et al., 2020)









Input





(Wang et al., 2020)



## **Definition of Loss Function**

• The loss of our network measures the cost incurred from incorrect predictions

 $L(f(x^{(i)}, W), y^{(i)})$ 

Predicted Actual

- Empirical Loss: Measures the total loss over our entire dataset
- Binary Cross Entropy Loss: Can be used with models that output a probability between 0 and 1
- Mean Squared Error Loss: Can e used with regression models that output continuous real numbers

### Training a Neural Network

- How can we use the loss function to train weights?
- We want to find the network weights that achieve the lowest loss
- We want to compute the lowest point and compute the gradient (iterative process)



# Computing Gradients: Backpropagation

• How much small change in w is going to affect our loss J



Repeat for every weight in the network using gradients from later layers

### **Training loss**

- Training a model means determining good values for all the weights and the bias from labeled examples
- □ Loss is the penalty for a bad prediction
- □ Loss is a number indicating how bad the model's prediction was on a single example
- □ If the model's prediction is perfect, the loss is zero, otherwise, the loss is greater
- □ Goal of training a model is to find a set of weights and biases that have loss, on average, across all examples



[Machine Learning Crash Course Google, 2020]









#### Step 1: Preprocessing



(Satya Mallick, 2016)

Step 2: Calculate the Gradient Images

$$g = \sqrt{g_x^2 + g_y^2}$$

$$\theta = \arctan \frac{g_y}{g_x}$$



Left : Absolute value of x-gradient. Center : Absolute value of ygradient. Right : Magnitude of gradient.

(Satya Mallick, 2016)

#### Step 3: Calculate Histogram of Gradients





(Satya Mallick, 2016)



(Satya Mallick, 2016)



(Satya Mallick, 2016)







(a) Industrial Buildings



(b) Residential Buildings







(h) Highway



(d) Permanent Crop



(i) Pasture



(e) River



(j) Forest

(Helber. P, et.al, 2015)



(f) Sea & Lake



(g) Herbaceous Vegetation



## 13 Bands covered by Sentinel's 2 Multispectral band

Band	Spatial Resolution m	<b>Central</b> Wavelength nm
B01 - Aerosols	60	443
B02 - Blue	10	490
B03 - Green	10	560
B04 - Red	10	665
B05 - Red edge 1	20	705
B06 - Red edge 2	20	740
B07 - Red edge 3	20	783
B08 - NIR	10	842
B08A - Red edge 4	20	865
B09 - Water vapor	60	945
B10 - Cirrus	60	1375
B11 - SWIR 1	20	1610
B12 - SWIR 2	20	2190



Model	Layers	Hidden size D	MLP size	Heads	Params
ViT-Base	12	768	3072	12	86M
ViT-Large	24	1024	4096	16	307M
ViT-Huge	32	1280	5120	16	632M

Table 1: Details of Vision Transformer model variants.

Model	Layers	Hidden size D	Heads	Image size (px)
ViT small	12	384	6	224×224





#### Normalization variation reconstruction



b)

c)

d)

#### Masking variation reconstruction

0.7

0.5

0.2












_	RAW norm	HOG norm+norm.hog	Sup. Learning
Annual Crop	99.0	98.83	97.16
Forest	99.5	99.83	99.67
Herbaceous Vegetation	98.16	98.67	98.16
Highway	97.4	97.6	94.6
Industrial	97.8	90.0	97.6
Pasture	99.0	98.0	96.0
Permanent Crop	98.2	97.4	97.0
Residential	99.5	99.5	98.67
River	99.0	98.8	98.0
Sea Lake	100.0	99.67	99.0

## Table 5.5 Confusion Matrix: Accuracies per class for 100% of the labels

-	RAW norm	HOG norm+norm.hog	Sup. Learning
Annual Crop	91.83	93.0	87.33
Forest	99.0	99.17	96.67
Herbaceous Vegetation	94.5	95.83	91.5
Highway	87.0	87.4	75.2
Industrial	96.8	98.0	93.4
Pasture	93.75	94.0	89.25
Permanent Crop	89.8	90.2	79.2
Residential	95.67	98.0	90.83
River	97.4	97.0	95.6
Sea Lake	95.17	97.33	96.17

## Table 5.6 Confusion Matrix: Accuracies per class for 10% of the labels



Figure 5.22 (a) Highway misclassified as permanent crop, (b) Highway misclasified as permanent crop, (c) Highway misclassified as Industrial zone



Figure 5.23 (a) Industrial zone misclassifeid as river, (b) Residential zone misclassified as permanent crop, (c) Industrial zone misclassified as highway



Figure 5.24 (a) Pasture misclassified as Herb. Vegetation, (b) Annual crop misclassified as permanent crop, (c) Herb. vegetation misclassified as permanent crop



Figure 5.25 (a) Permanent crop misclassified as highway, (b) Sea lake misclassified as river, (c) River misclassified as highway



(Alexander Amini, 2022)



(Ishan Misra, 2021)