#### Joint Energy-based Model for Remote Sensing Image Processing

Daniel Orozco Helmholtz Al Conference 2023 Hamburg, 13.06.2023



### Wissen für Morgen

#### Content

- 1. Introduction
- 2. Joint Energy-based Model (JEM)
- 3. JEM in Remote Sensing
  - 3.1 JEM instability for Classification
  - 3.2 Extension of JEM for Segmentation
- 4. Conclusions and future work



#### 1. Introduction

#### Earth Observation Data

#### – Volume

- ESA Sentinel satellites: ~20TiB daily user-level data.1
- Deep Learning models with millions of parameters.

#### - Quality

- Completely labeled sets are costly.
- Require a long time to update.
- High-resolution semantic segmentation maps.



**2. JEM** 



#### 2.1 Energy function

#### Joint Energy-based Model

- Energy function
  - ↑ High values: unrealistic samples.
  - $\downarrow$  Low values: realistic samples.

#### **Dual optimization problem**





#### 3. JEM in Remote Sensing

#### Datasets

#### - EuroSAT (RGB)

- 22,000 (train) and 5,000 (validation & test).
- 10 classes.



[7] Davies

#### - New York City (RGB + NIR)

- 6,000 (train) and 1,140 (validation & test).
- 3 classes.
  - Buildings, Vegetation and Bare land.

#### Optical



#### Ground Truth





3. JEM in Remote Sensing

# RQ1: What factors influence the <u>instability</u> of JEM training for Classification and how to alleviate it?



#### 3.1 JEM training instability

– The training loss  $(\mathcal{L})$  indicates how well the model fits the training data.

– JEM training simultaneously optimizes the classification and generation branch:

 $\mathcal{L}_{TOTAL} = \mathcal{L}_{CLA} + \mathcal{L}_{GEN}$ 

$$\mathcal{L}_{CLA} = CrossEntropy(y, \hat{y}) = -\sum_{k} y[k] \cdot \log(\hat{y}[k]) \implies \text{Bounded above Optimized}$$
$$\mathcal{L}_{GEN} = E_{\theta}(x) - E_{\theta}(x') \implies \text{Unbounded}$$

– Model **divergence** due to unbounded definition of generation loss.



#### 3.1 JEM training instability

- Model **divergence** due to unbounded definition of generation loss

Dataset	Completed epochs	Completed epochs (%)
CIFAR-10	51	25.5%
EuroSat	10	5%

Table 3.2. Completed epochs for JEM classification with different datasets.





3.1 Regularization

- **Regularization** adds a penalty term to the objective function.

– Applied to generation loss to **constrain** the values:

 $\mathcal{L}_{GEN\_REG} = \mathcal{L}_{GEN} + \alpha \cdot REG$ 

– Three types of **L2** regularization:

$$REG_r = \sum E_\theta(x)^2$$

(Energy values of training/real samples)

$$REG_u = \sum E_\theta(x')^2$$

 $REG_g = \sum \nabla E_{\theta}(x)(x)^2$ 

(Energy values of generated/unreal samples)



#### 3.1 Regularization

#### - Regularization results

Regularization			pleted oochs	Best classification	
Weight	Туре	#	%	accuracy (Validation)	
	Energy of training samples REG <sub>r</sub>	200	<b>100%</b>	95.16%	
0.2	Energy of generated samples REG <sub>u</sub>	155	77.5%	63.42%	
	$\nabla$ Energy of training samples $REG_g$	188	59%	89.82%	

Table 3.5. L2 regularization on energy values of training and generated samples and on energy gradients of training samples.

Regularization type	Enables more training epochs	Complete training	Average time per epoch (seconds)
Energy of training samples REG <sub>r</sub>		~	288
Energy of generated samples REG <sub>u</sub>	V	289	
$\nabla$ Energy of training samples $REG_g$		^	398



Table 3.6. Regularization types implications.

#### 3.1 Regularization

#### - Regularization and hyper-parameters results

Learning	Regularization		Completed epochs		Best classification
Tate	Type	Weight	#	%	accuracy (variation)
		0	10	5%	84%
$1 \times 10^{-4}$	Energy of training samples	0.5	7	3.5%	54%
		1	6	3%	40%
		0	86	43%	88.6%
	Energy of training samples	0.2	200	100%	95.16%
$1 \times 10^{-5}$		0.5	31	15.5%	66.8%
		0.75	30	15%	63%
	Energy of generated samples	0.2	155	77.5%	63.42%
	$\nabla$ Energy of training samples	0.2	118	59%	89.82%



Table 3.10. Completed epochs for different hyper-parameters and regularizations.

3. JEM in Remote Sensing

## RQ2: What are the potential and challenges of JEM for <u>semantic segmentation</u> of remote sensing data?





#### - Vanilla U-Net results

Architecture	Training size (samples)	Validation size (samples)	Sample size (pixels)	Best segmentation accuracy (validation)
	6,000	1,140	256x256	29.67%
Vanilla U-Net			128x128	55.94%
			64x64	56.62%

Table 4.2. Segmentation results with U-Net and different sample sizes.



#### - Wide-ResNet + U-Net results:

• +10% accuracy w.r.t. Vanilla U-Net architecture.

Architecture	Sample size (pixels)	Training size (samples)	Validation size (samples)	Best segmentation accuracy (validation)
Wide-ResNet		10	5	67.88%
+	128x128	80	20	62.43%
U-Net		2,000	20	47.50%

Table 4.8. Segmentation results with Wide-ResNet + U-Net and different sizes for training and validation sets.



- Wide-ResNet + U-Net results





Epoch 6

Epoch 3





Epoch 1

- Divergence

• Pixel aggregations produce energy values ~4 orders of magnitude bigger w.r.t. classification.





#### 4. Conclusions and future work

#### Divergence:

- Regularization stabilizes the training and alleviates the divergence to some extent.
- Regularization eases the competition of the dual optimization problem.

#### sJEM:

- The pixel independence assumption is too strong for RS images where spatial autocorrelation is a common characteristic.
- Stronger regularization techniques are required for RS image segmentation.

#### Future work:

- Further regularization terms that are more effective for divergence.
- Robust definitions to aggregate pixels' energy values, including measures of spatial autocorrelation.
- Feedback sensibility in the generation process to "bad-looking" samples.



Earth Observation Center

## Thanks for your attention!

## Questions?



#### References

[1] W. Grathwohl, K.-C. Wang, J.-H. Jacobsen, D. Duvenaud, M. Norouzi, and K. Swersky, «Your classifier is secretly an energy-based model and you should treat it like one », Proceedings of the International Conference on Learning Representations (ICLR), 2020.

[2] S. Zhao, J.-H. Jacobsen, and W. Grathwohl, « Joint energy-based models for semi-supervised classification », in Proceedings of the International Conference on Machine Learning - Workshop on Uncertainty and Robustness in Deep Learning (ICMLW), 2020.

[3] J. Castillo-Navarro, B. Le Saux, A. Boulch, and S. Lefèvre, « Energy-based models in Earth observation: from generation to semi-supervised learning », IEEE Transactions on Geoscience and Remote Sensing, 2021.



#### References

[4] Serco Europe, RP. (2021). Sentinel data access annual report Y2020. Copernicus. <u>https://scihub.copernicus.eu/twiki/pub/SciHubWebPortal/AnnualReport2020/COPE-SERCO-RP-21-1141\_-</u> <u>Sentinel\_Data\_Access\_Annual\_Report\_Y2020\_final\_v2.3.pdf</u>

[5] P. Helber, B. Bischke, A. Dengel, and D. Borth. "Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification." In: IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing 12.7 (2019), pp. 2217–2226.

[6] C. M. Albrecht, F. Marianno, and L. J. Klein. "Autogeolabel: Automated label generation for geospatial machine learning." In: 2021 IEEE International Conference on Big Data (Big Data). IEEE. 2021, pp. 1779–1786.

[7] Davies. A, « Computer Vision | Land Cover Classification Using TensorFlow in Python », in Towards Data Science, 2022.



- **Regularization** in the view of energy landscape



– Influence of **hyperparameters**: learning rate

Learning rate	Completed epochs	Completed epochs (%)	Best classification accuracy (validation)
$1 \times 10^{-4}$	10	5%	84%
$1 \times 10^{-5}$	86	43%	88.6%

Table 3.8. Completed epochs for different learning rates.



– Influence of **hyperparameters**: SGLD steps

Number of SGLD steps	Average time per epoch (seconds)	Completed epochs	Completed epochs (%)
40	2310	10	5%
100	5080	9	4.5%

Table 3.9. Completed epochs and average time per epoch for different number of SGLD steps.



Figure 3.11. Generated samples for different number of SGLD steps.



- Detection and Restart





#### - Detection and Restart





Table 3.3. Thresholds for divergence detection.

#### - Detection and Restart



#### - Detection and Restart

Algorithm 1 Pre-emptive	e detection of divergence and restart of training
$\alpha_{RSD} \leftarrow 0.3$	
$\alpha_{OTL} \leftarrow 15$	
$jump\_size \leftarrow 30$	
$\omega \leftarrow 10$	window size for the rolling standard deviation
$current\_rstd \leftarrow Rolling$	$gStd( x-x' ,\omega)$
$num\_outliers \leftarrow \sum_x IQ$	$R_Outlier(x)$
<b>if</b> ( <i>current_rstd</i> $\geq \alpha_{RSI}$	b) & $(num\_outliers \ge \alpha_{OTL})$ then
Divergence warnin	g!
previous_replay_bu	$ffer \leftarrow LoadBuffer(jump_size)$
previous_optimizer	$\leftarrow$ LoadOptimizer(jump_size)
$previous\_model \leftarrow$	LoadModel(jump_size)
RestartTraining(pro	evious_replay_buffer, previous_optimizer, previous_model)
end if	



#### Appendix | JEM for Classification

#### – Potential for **semi-supervised learning**



#### Appendix | JEM for Classification

#### – Potential for **semi-supervised learning**

Mode (% labels)	Configuration	Mean accuracy	ECE
Supervised (100% labels)	JEM	95.16%	1.08%
Supervised (100% labels)	Classification	94.86%	3.37%
Comi and (200/ labela)	JEM	94.06%	0.71%
Semi-supervised (80% labels)	Classification	94.42%	4.10%
Comi our orrigod (200/ labola)	JEM	80.76%	0.91%
Semi-supervised (20% labels)	Classification	81.02%	15.05%

Table 3.13. EuroSat supervised vs semi-supervised classification accuracy and ECE.



#### Appendix | JEM for Segmentation (sJEM)

- Individual **segmentation** branch analysis:
  - U-Net architecture.
  - Balance between training time and accuracy.
  - Dice metric addition to segmentation loss.

Sample size (pixels)	Average time per epoch (seconds)	Best segmentation accuracy (validation)
64x64	11	83.83%
128x128	21	86.91%
256x256	54	88.67%

Table 4.3. Segmentation branch results with Vanilla U-Net and different sample sizes.

		Sample size	
_	64x64	128x128	256x256
Input image			
Ground Truth Map			
– Predicted Map			
	Buildings	Vegetation	Bare land

Sample size (pixels)	Segmentation Loss	Best segmentation accuracy (validation)
64x64	Cross-Entropy	83.83%
	Cross-Entropy + Dice	84.40%
128x128	Cross-Entropy	86.91%
	Cross-Entropy + Dice	87.05%
256x256	Cross-Entropy	88.67%
	Cross-Entropy + Dice	88.72%

Table 4.4. Segmentation branch results with different loss metrics.



#### Appendix | JEM for Segmentation (sJEM)

– Individual **generation** branch analysis:

- 1 sample.
- Different encoders and learning rates.





#### Appendix | JEM for Segmentation (sJEM)

– Influence of hyperparameters: learning rate and buffer size

Learning rate	Buffer size	Best segmentation accuracy (validation)
$1 \times 10^{-4}$	3,000	67.88%
	6,000	59.48%
	9,000	39.25%
	30,000	38.26%
	60,000	38.37%
$1 \times 10^{-5}$	3,000	39.02%
	6,000	38.86%

Table 4.7. Segmentation results with Wide-ResNet + U-Net, different learning rates, and buffer sizes.

