

Afro-MNIST

Wu, Yang,
Prabhu

PML4DC
ICLR 2020

Introduction

Proposal

Contributions

Methodology

Experiments

Summary

References

Afro-MNIST

Synthetic generation of MNIST-style datasets for
low-resource languages

Daniel J. Wu¹ **Andrew C. Yang**¹ Vinay Prabhu²

¹Stanford University

²UnifyID Inc.

Introduction

Afro-MNIST

Wu, Yang,
Prabhu

PML4DC
ICLR 2020

Introduction

Proposal

Contributions

Methodology

Experiments

Summary

References

- Classifying Hindu-Arabic numerals in the MNIST dataset¹ has become the “Hello world” challenge in the machine learning community.
- This task has excited a large number of prospective machine learning scientists and has led to practical advancements in OCR.



0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9 9 9

MNIST dataset²

¹Yann LeCun et al. “Gradient-based learning applied to document recognition”. In: *Proceedings of the IEEE* 86.11 (1998), pp. 2278–2324.

²Wikimedia Commons.

Introduction

Afro-MNIST

Wu, Yang,
Prabhu

PML4DC
ICLR 2020

Introduction

Proposal

Contributions

Methodology

Experiments

Summary

References

- Work in ML and AI focuses almost exclusively on high-resource languages, which use the Hindu-Arabic numeral system.
- Of over 7,000 languages in the world,³ the vast majority are not represented in the ML research community.
- There are many alternative numeral systems for which an MNIST-style dataset is not available.

³David M. Eberhard, Gary F Simons, and Charles D Fennig. *Languages of the World*. 2019. URL: <http://www.ethnologue.com/>.

Introduction

Afro-MNIST

Wu, Yang,
Prabhu

PML4DC
ICLR 2020

Introduction

Proposal

Contributions

Methodology

Experiments

Summary

References

- Much of the world's linguistic diversity comes from languages spoken in developing nations.
- There is a wealth of linguistic diversity in the languages of Africa, many of which have dedicated orthographies and numeral systems.
- One notable example is the Ge'ez (Ethiopic) script, which is used to transcribe languages such as Amharic and Tigrinya, spoken by some 30 million people.⁴

ገዕዝ

Ge'ez written in the Ge'ez script

⁴*Ibid.*

Introduction

Afro-MNIST

Wu, Yang,
Prabhu

PML4DC
ICLR 2020

Introduction

Proposal

Contributions

Methodology

Experiments

Summary

References

The numeral system is the most endangered aspect of any language⁵

⁵Emmanuel Mfanafuthi Mgqwashu. "Academic literacy in the mother tongue: A pre-requisite for epistemological access". In: *Diversity, Transformation and Student Experience in Higher Education Teaching and Learning* (2011), p. 159.

Proposal

Afro-MNIST

Wu, Yang,
Prabhu

PML4DC
ICLR 2020

Introduction

Proposal

Contributions

Methodology

Experiments

Summary

References

- Large amounts of training data for African languages such as these are not readily available.
- But effective neural networks can be trained on highly perturbed versions of just a single image of each class.⁶
- We experiment with creating synthetic numerals that mimic the likeness of hand-written numerals in those writing systems.

⁶Alexey Dosovitskiy et al. "Discriminative unsupervised feature learning with exemplar convolutional neural networks". In: *IEEE transactions on pattern analysis and machine intelligence* 38.9 (2015), pp. 1734–1747.

Contributions

Afro-MNIST

Wu, Yang,
Prabhu

PML4DC
ICLR 2020

Introduction

Proposal

Contributions

Methodology

Experiments

Summary

References

- We release synthetic MNIST-style datasets for four scripts used to write Afro-Asiatic or Niger-Congo languages: Ge'ez, Vai, Osmanya, N'Ko⁷, which serve as **drop-in replacements** for the MNIST dataset.
- We describe a general framework for resource-light syntheses of MNIST-style datasets.
- These datasets can be found at <https://github.com/daniel-wu/AfroMNIST>.

⁷The Vai, Osmanya, and N'Ko scripts are not in wide use, but nonetheless they can be synthesized using the methods we present.

Methodology

Afro-MNIST

Wu, Yang,
Prabhu

PML4DC
ICLR 2020

Introduction

Proposal

Contributions

Methodology

Experiments

Summary

References

- Generate an exemplar seed dataset for each numeral system from the corresponding Unicode characters.
- Apply series of elastic deformations and corruptions.⁸



⁸We note that, in cases where a limited amount of handwritten data is available, deformations and corruptions can be applied to those examples instead of Unicode exemplars.

Experiments

Afro-MNIST

Wu, Yang,
Prabhu

PML4DC
ICLR 2020

Introduction

Proposal

Contributions

Methodology

Experiments

Summary

References

- We begin by training LeNet-5, the network architecture first used on the original MNIST dataset⁹, for numeral classification.

Dataset	Accuracy (%)
MNIST	99.65
Ge'ez-MNIST	99.92
Vai-MNIST	100
Osmanya-MNIST	99.99
N'Ko-MNIST	100

- After testing a LeNet-5 trained on Ge'ez-MNIST on a small dataset of handwritten Ge'ez numerals¹⁰, we found the model achieved only an accuracy of 30.30%.

⁹LeCun et al., "Gradient-based learning applied to document recognition".

¹⁰Tesfamichael Molla. *Ethiopian-MNIST*.

<https://github.com/Tesfamichael1074/Ethiopian-MNIST>. 2019.

Experiments

Afro-MNIST

Wu, Yang,
Prabhu

PML4DC
ICLR 2020

Introduction

Proposal

Contributions

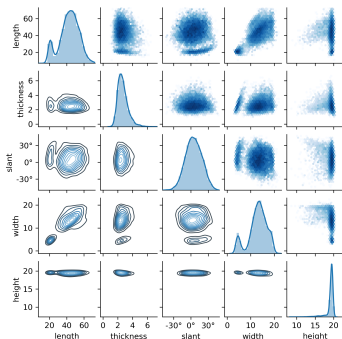
Methodology

Experiments

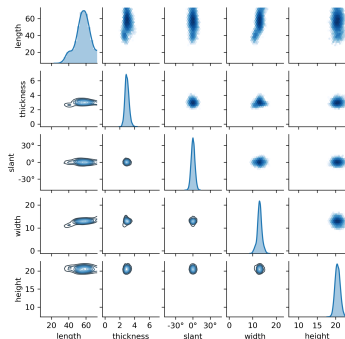
Summary

References

- Morphological comparisons¹¹ between MNIST and Ge'ez-MNIST show clear differences in variance.



MNIST



Ge'ez-MNIST

¹¹Vinay Uday Prabhu. "Kannada-mnist: A new handwritten digits dataset for the kannada language". In: *arXiv preprint arXiv:1908.01242* (2019).

Summary

Afro-MNIST

Wu, Yang,
Prabhu

PML4DC
ICLR 2020

Introduction

Proposal

Contributions

Methodology

Experiments

Summary

References

- Many writing systems, especially those used in developing nations, are underrepresented in the ML community.
- Elastic deformations and corruptions show promise in generating synthetic numeral data, but other methods of creating synthetic digits more similar to handwritten digits ought to be explored as well.
- We expect this benchmark to be a fertile starting point for exploring augmentation and transfer learning strategies for low-resource languages.
- We hope that endeavors such as these help encourage the next generation of diverse ML practitioners to be part of the broader research community.

References

Afro-MNIST

Wu, Yang,
Prabhu

PML4DC
ICLR 2020

Introduction

Proposal

Contributions

Methodology

Experiments

Summary

References

Dosovitskiy, Alexey et al. “Discriminative unsupervised feature learning with exemplar convolutional neural networks”. In: *IEEE transactions on pattern analysis and machine intelligence* 38.9 (2015), pp. 1734–1747.

Eberhard, David M., Gary F Simons, and Charles D Fennig. *Languages of the World*. 2019. URL: <http://www.ethnologue.com/>.

LeCun, Yann et al. “Gradient-based learning applied to document recognition”. In: *Proceedings of the IEEE* 86.11 (1998), pp. 2278–2324.

Mgqwashu, Emmanuel Mfanafuthi. “Academic literacy in the mother tongue: A pre-requisite for epistemological access”. In: *Diversity, Transformation and Student Experience in Higher Education Teaching and Learning* (2011), p. 159.

Molla, Tesfamichael. *Ethiopian-MNIST*.
<https://github.com/Tesfamichael1074/Ethiopian-MNIST>. 2019.

Prabhu, Vinay Uday. “Kannada-mnist: A new handwritten digits dataset for the kannada language”. In: *arXiv preprint arXiv:1908.01242* (2019).