CNN FOR SPECTROSCOPIC REDSHIFT ESTIMATION ON EUCLID DATA

FORTH Institute of Computer Science





EUCLID SEDs



EUCLID SEDs



PREDICTIVE MODEL

Keusiiii

Real-valued, non-negative number (z)

Regression Analysis

- not very robust algorithms
- not state-of-the-art

Split the examined redshift interval into ordinal classes, based on Euclid's characteristic resolution

Classification Problem

- utilization of state-of-the-art classifiers, that fully harness the power of Big Data
- probabilistic approach using Softmax quantification of the levels of confidence for each estimation

DATASET & METHODOLOGY

- 10K initial clean rest-frame spectral profiles (SEDs)
- <u>Generation of randomly redshifted examples</u> $\circ z = [1, 1.8)$, similar to Euclid's specification $\circ \log(1 + z) = \log(\lambda_{observed}) - \log(\lambda_{emit}) \iff 1 + z = \frac{\lambda_{observed}}{\lambda_{emit}}$ $\circ \text{ optional addition of white Gaussian noise (idealistic vs. realistic)}$
- Quantization of the utilized redshift range

○ split into 800 discrete classes, which implies a resolution of 0.001

R. Stivaktakis, G. Tsagkatakis, B. Moraes, F. Abdalla, J-L Starck, P. Tsakalides, "Convolutional Neural Networks for Spectroscopic Redshift Estimation on Euclid Data", IEEE Transactions on Big Data (Special Issue on Big Data From Space, 2020)

CONVOLUTIONAL NEURAL NETWORKS

- inspired by the concept of "visual receptive fields"
- automated feature extractors
- exhibit spatial correlations of the given input

• less prone to overfitting

local-connectivity property & weight-sharing lead to a dramatic decrease in total parameters
effective utilization of Big Data results in a high-generalization capacity

• Dropout & Batch Normalization

1-DIMENSIONAL CNN

EXPERIMENTAL SETUP

Sample Size

- 400,000 training samples
- 10,000 CV samples
- 10,000 testing samples

Compared Classifiers

- k Nearest Neighbors
- Random Forests
- Support Vector Machines

Hyper-parameters

- 3 Conv. + ReLU Layers
- 1 Dense + Softmax Layer
- Kernel size = 8, Stride = 1
- # Filters = 16 per Conv. Layer
- Categorical Cross-Entropy Loss

SOME IDEAS....

- Same input -> different outputs
- Same outputs <- different inputs
- Different model architectures (ML, Bayesian, etc.)
- Noise in labels and/or input
- Uncertainty labels and/or input
- Classification vs regression
- Loss functions

• Few-shot learning / Imbalanced labels