# Galaxy Classification Challenge

**Group #1**

**Timo Schellhaas**
**Jaroslav Storek**
**Leonid Lunin**
**Lisa Lou Krümmel**
**Mathis Nolte**

# Dataset Exploration and the Class 5 Problem



Example images of each class from Galaxy10 dataset

| Disk, Face-on, No Spiral | Smooth, Completely round | Smooth, in-between round | Smooth, Cigar shaped | Disk, Edge-on, Rounded Bulge |
| Disk, Edge-on, Boxy Bulge | Disk, Edge-on, No Bulge | Disk, Face-on, Tight Spiral | Disk, Face-on, Medium Spiral | Disk, Face-on, Loose Spiral |

Galaxy10 Dataset: Henry Leung/Jo Bovy 2018, Data Source: SDSS/Galaxy Zoo

# Class 5 Problem

- Class 5 is dramatically underrepresented



Distribution of categories

| Category | Count |
|---|---|
| Disk, Face-on, No Spiral | 3461 |
| Smooth, Completely round | 6997 |
| Smooth, in-between round | 6292 |
| Smooth, Cigar shaped | 394 |
| Disk, Edge-on, Rounded Bulge | 1534 |
| Disk, Edge-on, Boxy Bulge | 17 |
| Disk, Edge-on, No Bulge | 589 |
| Disk, Face-on, Tight Spiral | 1121 |
| Disk, Face-on, Medium Spiral | 906 |
| Disk, Face-on, Loose Spiral | 519 |

# Our solutions

- Replicating class 5 images to increase its impact
- Data augmentation will prevent overfitting (random rotation and flip of images but no stretching)
- Increase batch size to have a higher possibility to actually see a Class 5

Distribution of categories

| Category | Count |
| --- | --- |
| Disk, Face-on, No Spiral | 3461 |
| Smooth, Completely round | 6997 |
| Smooth, in-between round | 6292 |
| Smooth, Cigar shaped | 394 |
| Disk, Edge-on, Rounded Bulge | 1534 |
| Disk, Edge-on, Boxy Bulge | 1020 |
| Disk, Edge-on, No Bulge | 589 |
| Disk, Face-on, Tight Spiral | 1121 |
| Disk, Face-on, Medium Spiral | 906 |
| Disk, Face-on, Loose Spiral | 519 |

# Neural Network Architecture

# Which network have we used and how do we measure the performance?

**Average accuracy among the classes** -to rate the overall performance

**Confusion matrix** - to rate the performance of the individual classes

**What we have tried:**

- LeNet
- AlexNet
- Branching

**Other possible approaches:**

- ResNet [1]
- Assembles

[1] Ba Alawi, A. E., & Al-Roainy, A. A. (2021). Deep Residual Networks Model for Star-Galaxy Classification. 2021 https://doi.org/10.1109/icoten52080.2021.9493433

# LeNet and AlexNet inspired architecture

- Dropout layers to reduce overfitting
- Optimizer: Adam
- Using L2 weight_decay option
- Kernel sizes matched to the feature sizes in the images
- Convolution kernel sizes: 5 and 7

```
----------------------------------------------------------------------
       Layer (type)              Output Shape           Param #
======================================================================
         Conv2d-1           [2048, 6, 64, 64]               156
           ReLU-2           [2048, 6, 64, 64]                 0
      AvgPool2d-3           [2048, 6, 32, 32]                 0
         Conv2d-4          [2048, 16, 28, 28]             2,416
           ReLU-5          [2048, 16, 28, 28]                 0
      AvgPool2d-6          [2048, 16, 14, 14]                 0
        Flatten-7                [2048, 3136]                 0
         Linear-8                 [2048, 120]           376,440
           ReLU-9                 [2048, 120]                 0
        Linear-10                  [2048, 84]            10,164
          ReLU-11                  [2048, 84]                 0
        Linear-12                  [2048, 10]               850
======================================================================
Total params: 390,026
Trainable params: 390,026
Non-trainable params: 0
----------------------------------------------------------------------
Input size (MB): 32.00
Forward/backward pass size (MB): 1360.53
Params size (MB): 1.49
Estimated Total Size (MB): 1394.02
----------------------------------------------------------------------
```
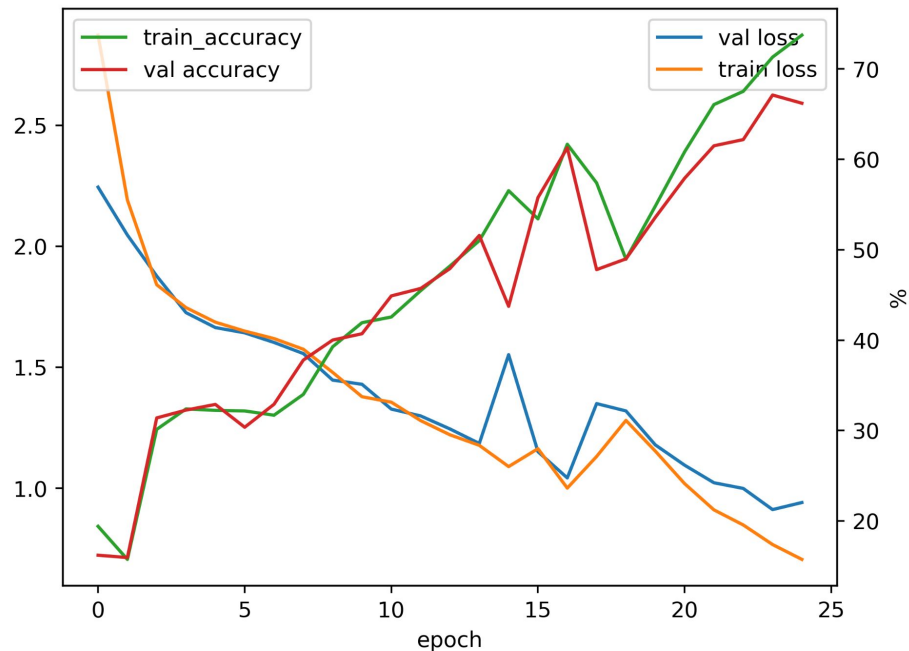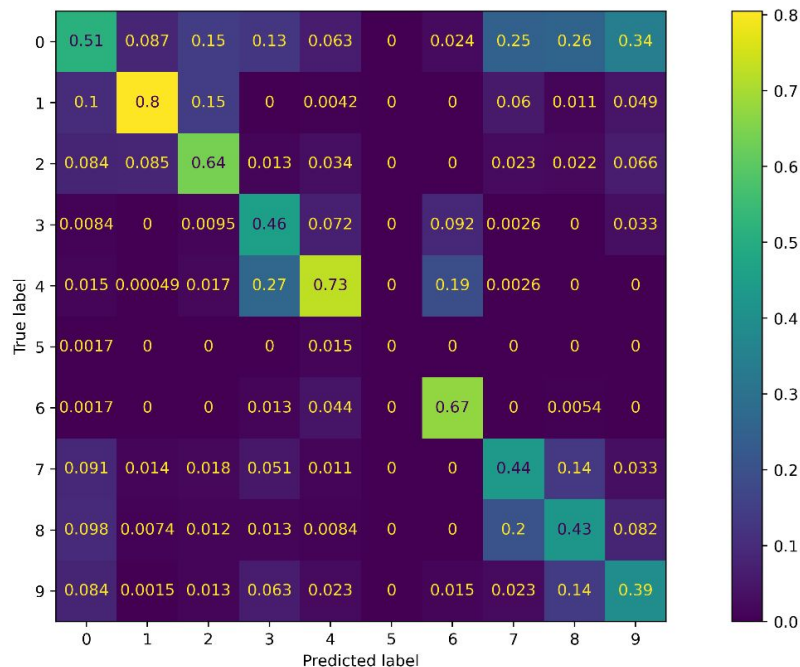
Network Training and Evaluation

# No data augmentation (200 epochs, 2048 batch size)

- Overfitting!

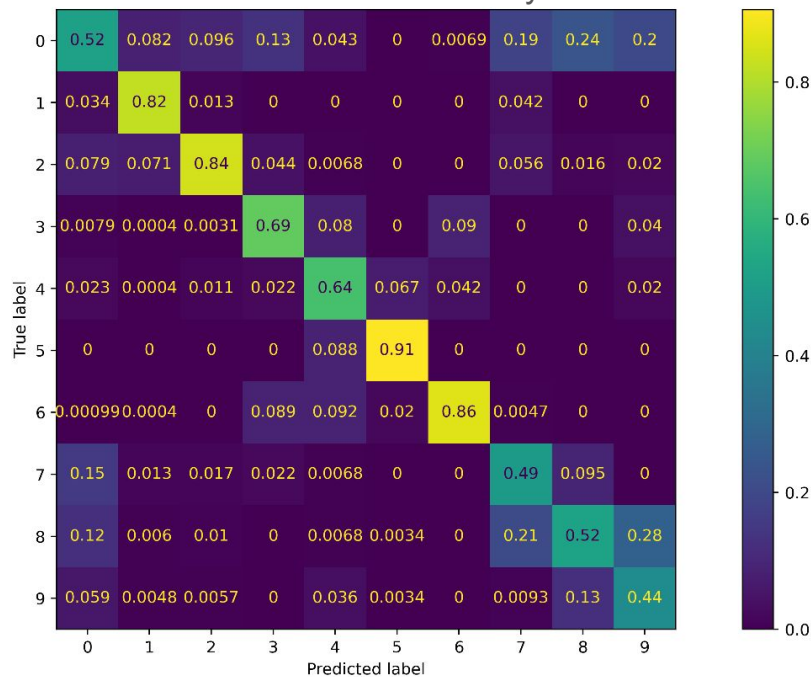# With data augmentation (200 epochs, 2048 batch size)

- Best result: 78% accuracy on the test dataset with 78% accuracy on train dataset
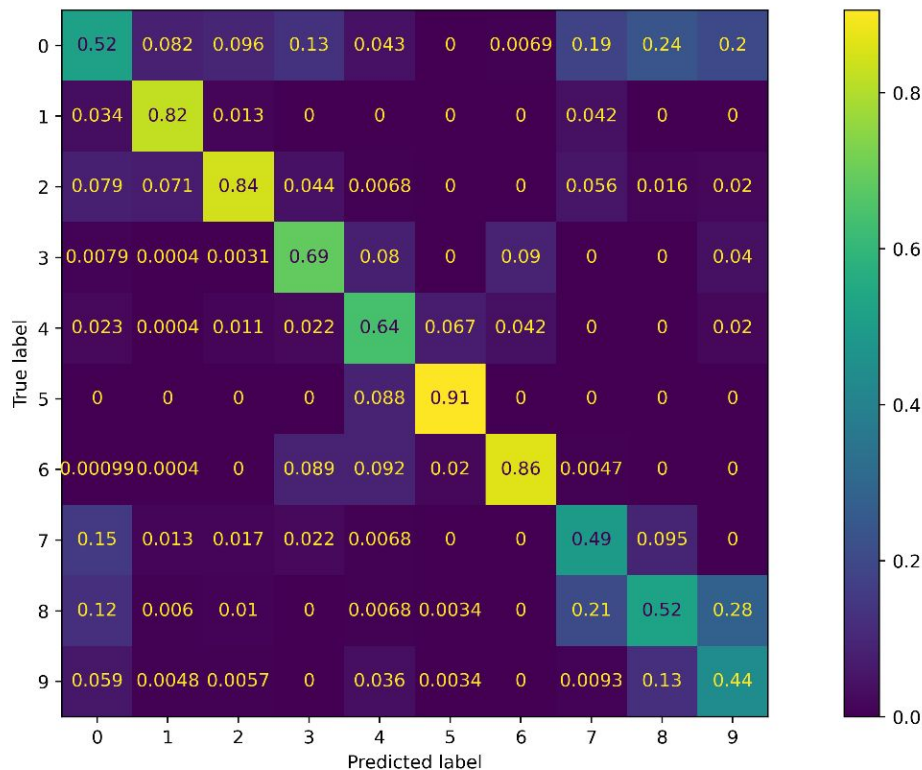
# Performance on the individual classes

- Excellent performance on classes with high statistics (1 & 2)
- Outstanding performance on class 5
- Worst performance on poorly represented classes with high similarity

**Potential for improvement:**

- Replicating the last three classes

# Thank you for your attention!

Jupyter notebook available at:

github.com/lrlunin/erum-datahub-challange-2023/tree/main