NNDL@UNIPD

# ENSEMBLING CNN MODELS FOR 3D OBJECT CLASSIFICATION BASED ON VOXEL GRID REPRESENTATION

## Authors

D'Antimo Simone[1]     simone.dantimo@studenti.unipd.it
Singh Harjot[1]        harjot.singh@studenti.unipd.it

[1]The University of Padua

Ensembling CNN models for 3D Object Classification based on Voxel Grid representation

Simone D'Antimo[†], Harjot Singh[†], Jacopo Pegoraro[†]

*Abstract*—Convolutional Neural Networks (CNN) have been widely used as Deep learning tools to enhance the spatial relationship among images and videos. In this paper, we illustrate the potential of CNN and Voxel grids for the 3D Object classification task. Intuitively, given a 3D image, we split it into blocks, each of those blocks may be only full or empty, which implies that the resolution of the 3D object is strongly correlated with the number of blocks in which they are split in, as trade-off more blocks implies more memory required for storage. We decided to evaluate the potentiality and the limits of a 3D Object Classification System having a limited amount of computational capacity and memory storage: those limitations imply that we must use samples with a lower resolution (i.e. composed by a limited number of voxel), loosing some details and shape of the objects, but at the same time the results obtained show that we created a robust model which perform efficient analysis and can classify the objects independently of their rotation. We were able to reach good results on the test set, with an accuracy rate of 88 percent at the end.

*Index Terms*—Neural Networks, Convolutional Neural Networks, 3D Object Classification, Voxel Grid, ModelNet10, Data Augmentation.

# TABLE OF CONTENTS

THE **CHALLENGE**

# THE CHALLENGE

**Our goal:** Classifying 3D Object independently their rotation using CNN and Voxel Grid representation

**Our data:** ModelNet10,
- 4899 3D CAD divided into 10 classes
- 80/20* = train/test* splitting
- orientation aligned

- CAD models
- Non i.i.d.
- Desired partition 80-10-10

## TRAIN SET

■ Initial samples



| | bathtub | bed | chair | desk | dresser | monitor | night_stand | sofa | table | toilet |
|---|---|---|---|---|---|---|---|---|---|---|
| ■ Initial samples | 106 | 515 | 889 | 200 | 200 | 465 | 200 | 680 | 392 | 344 |

# OUR APPROACH

Mesh

Voxelizzation

Data Augmentation

Partitioning

Weak model

ENSEMBLE

Prediction

Data preparation

Model

# DATA PREPARATION

# DATA PREPARATION

## VOXELIZATION

**1** ModelNet10
3D Cad Models



Voxel Grids

# DATA PREPARATION

**1** **Focus** on making the **dataset IID**

**2** **N Rotations,** depending on their initial distribution

| Rotations | Classes |
|-----------|---------|
| 2 | bed, chair, monitor, sofa |
| 3 | table, toilet |
| 5 | desk, dresser, night_stand |
| 8 | bathtub |

# DATA PREPARATION

VOXELIZATION   >   **AUGMENTATION**   >

**1** **Focus** on making the **dataset IID**

**2** **N Rotations,** depending on their initial distribution

| Rotations | Classes |
|-----------|---------|
| 2 | bed, chair, monitor, sofa |
| 3 | table, toilet |
| 5 | desk, dresser, night_stand |
| 8 | bathtub |

## TRAIN SET

■ AFTER AUGMENTING   ■ Initial samples

| | bathtub | bed | chair | desk | dresser | monitor | night_stand | sofa | table | toilet |
|---|---------|-----|-------|------|---------|---------|-------------|------|-------|--------|
| ■ AFTER AUGMENTING | 848 | 1000 | 1000 | 1000 | 1000 | 930 | 1000 | 1000 | 1176 | 1032 |
| ■ Initial samples | 106 | 515 | 889 | 200 | 200 | 465 | 200 | 680 | 392 | 344 |

# DATA PREPARATION

**1** **Focus** on making the **dataset ready for the ENSEMBLE models**

**2** Each weak model should learn
- well on their part of the dataset
- some information on data of other models.

TABLE 2: Our Bagging datasets compositions

| Dataset | Chunk's subsets | Copied Subset |
|---|---|---|
| 1 | 01, 02, 03, 04, 05 | 10, 15, 20, 25 — 24 |
| 2 | 06, 07, 08, 09, 10 | 01, 11, 16, 21 — 05 |
| 3 | 11, 12, 13, 14, 15 | 02, 07, 17, 22 — 06 |
| 4 | 16, 17, 18, 19, 20 | 03, 08, 13, 23 — 12 |
| 5 | 21, 22, 23, 24, 25 | 04, 09, 14, 19 — 18 |

Weak model

Partitioning

ENSEMBLE

Prediction

Model

# DATA PREPARATION

VOXELIZATION   >   AUGMENTATION   >   PARTITIONING   >   **CUTOUT**

**1** **Masking Out** a random portion of the input data

- Improve robustness
- Improve generalization



Cutout

Reference: *Research Gate*

MODEL

# THE WEAK MODELS STRUCTURE



| Layer (type) | Output Shape | Param # |
|---|---|---|
| Conv3d-1 | [-1, 32, 32, 32, 32] | 896 |
| BatchNorm3d-2 | [-1, 32, 32, 32, 32] | 64 |
| ReLU-3 | [-1, 32, 32, 32, 32] | 0 |
| Dropout-4 | [-1, 32, 32, 32, 32] | 0 |
| MaxPool3d-5 | [-1, 32, 16, 16, 16] | 0 |
| Conv3d-6 | [-1, 64, 16, 16, 16] | 55,360 |
| BatchNorm3d-7 | [-1, 64, 16, 16, 16] | 128 |
| ReLU-8 | [-1, 64, 16, 16, 16] | 0 |
| Dropout-9 | [-1, 64, 16, 16, 16] | 0 |
| MaxPool3d-10 | [-1, 64, 8, 8, 8] | 0 |
| Conv3d-11 | [-1, 128, 8, 8, 8] | 221,312 |
| BatchNorm3d-12 | [-1, 128, 8, 8, 8] | 256 |
| ReLU-13 | [-1, 128, 8, 8, 8] | 0 |
| Dropout-14 | [-1, 128, 8, 8, 8] | 0 |
| MaxPool3d-15 | [-1, 128, 4, 4, 4] | 0 |
| Linear-16 | [-1, 128] | 1,048,704 |
| ReLU-17 | [-1, 128] | 0 |
| Dropout-18 | [-1, 128] | 0 |
| Linear-19 | [-1, 10] | 1,290 |

!! Represent in 3D for model flow, but 3DConv layer has 4 dimension

# TRAINING & EXPERIMENTAL RESULTS

# TRAINING WEAK MODELS

Epochs per model:100

- Converge quickly <50 Epochs

Learning rate: 1e-3

Optimizer: ADAM

Loss: Cross Entropy

Max Pooling Kernel: 2x2x2

Max Pooling Stride: 2x2x2

# RESULTS ON VALIDATION SET

| CLASS | M0 | M1 | M2 | M3 | M4 | EN |
|---|---|---|---|---|---|---|
| 0 bathtub | 0.98 | 0.98 | 0.94 | 0.97 | 0.93 | 0.98 |
| 1 bed | 0.92 | 0.90 | 0.92 | 0.96 | 0.94 | 0.94 |
| 2 chair | 0.96 | 0.91 | 0.98 | 0.94 | 0.96 | 0.97 |
| 3 desk | 0.62 | 0.67 | **0.80** | 0.75 | 0.65 | 0.78 |
| 4 dresser | 0.74 | 0.73 | 0.71 | 0.83 | 0.74 | 0.77 |
| 5 monitor | 0.94 | 0.97 | 0.97 | 0.96 | 0.97 | 0.96 |
| 6 night_stand | 0.71 | 0.75 | 0.75 | 0.81 | 0.67 | 0.75 |
| 7 sofa | 0.88 | 0.76 | 0.81 | 0.87 | 0.90 | 0.88 |
| 8 table | 0.67 | 0.77 | 0.81 | 0.87 | 0.90 | 0.88 |
| 9 toilet | **0.99** | 0.90 | 0.95 | 0.91 | 0.98 | 0.97 |

# RESULTS ON TEST SET

**Ensembled model results:**

- Good diagonal Confusion matrix
- Errors are "Symmetric"



*Confusion matrix of the ensemble model*

# RESULTS ON TEST SET

| | | ACCURACY |
|---|---|---|
| 1 | M0 | 0.83 |
| 2 | M1 | 0.83 |
| 3 | M2 | 0.85 |
| 4 | M3 | 0.87 |
| 5 | M4 | 0.86 |
| 6 | EN | **0.88** |

BASELINE ACCURACY: 0.90

ShapeNet: 0.83
DeepPano 0.85
VoxNet 0.92
Orion 0.938

19

# CONCLUSIONS AND FUTURE WORK

# CONCLUSIONS

Model has good accuracy

Dynamic, can be adapted on the computer in which is run

# FUTURE WORK

Adding more data and Models

Implements a way to capture long-range context information

END THANK YOU

NNDL_Project.ipynb

File  Modifica  Visualizza  Inserisci  Runtime  Strumenti  Guida

Commenta   Condividi

+ Codice   + Testo      Connetti

# Imports and Global Vars

Import all the packages

Import all the packages

You need pillow to plot, and MatplotLib requires the latest version therefore you need to update it

If you need the Pillow, Don't panic the fist time the env is going to be stopped purposefully to update the Pillow package for the plots.

*You need to re-run this cell manually if NEED_PILLOW = True.

NEED_PILLOW: True

Mostra codice

Installing Open3d
420.5/420.5 MB 3.7 MB/s eta 0:00:00
10.3/10.3 MB 101.6 MB/s eta 0:00:00
77.1/77.1 kB 8.6 MB/s eta 0:00:00
138.3/138.3 kB 15.5 MB/s eta 0:00:00
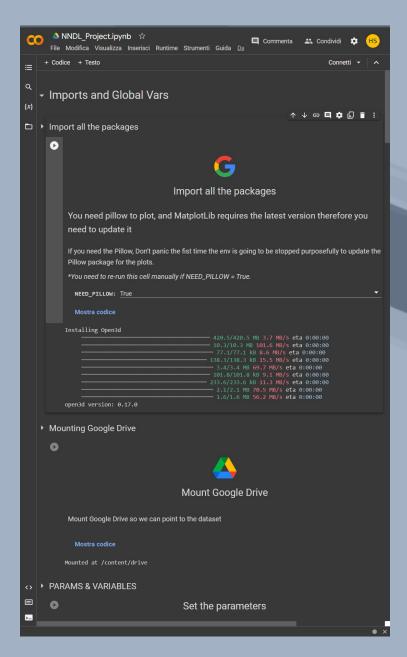3.4/3.4 MB 69.7 MB/s eta 0:00:00
101.8/101.8 kB 9.1 MB/s eta 0:00:00
233.6/233.6 kB 11.3 MB/s eta 0:00:00
2.1/2.1 MB 70.5 MB/s eta 0:00:00
1.6/1.6 MB 56.2 MB/s eta 0:00:00
open3d version: 0.17.0

Mounting Google Drive

Mount Google Drive

Mount Google Drive so we can point to the dataset

Mostra codice

Mounted at /content/drive

PARAMS & VARIABLES

Set the parameters

**NOTEBOOK DEMO**

NNDL_Project_BAGGING_kag...   Draft saved

Share   Save Version   9

File  Edit  View  Run  Add-ons  Help

+  🗑  ✂  ⧉  📋   ▶  ⏭ Run All   Markdown ▾   ● Draft Session off (run a cell to start)

Parameters/Variables for the Training Session

+ Code   + Markdown

```python
# Model need to be trained
MODEL = "VGNet_v3" #@param ["VGNetOrion_extended", "VGNet_v3"]{type:"string"
OPTIMIZER = "Adam" # @param ["Adam", "SGD"]{type:"string"}
LEARNING_RATE = 1e-3 # @param ["1e-2","1e-3","1e-4"]{type:"raw"}
#@markdown <pre>Normally LR would be:
#@markdown          <b>0.001</b> for ADAM
#@markdown          and <b>0.0001</b> with SGD</pre>

# Put to false if you need to do only the evaluation
TRAIN = False
BAGGING = True # for bagging (TODO: SEPARATE EPOCHS & old)
MODELS_TO_TRAIN = [0,1,2,3,4] # max count 2 (only 2 gpus available) could be

# If FALSE: 1) disables tqdm LOGGING
#           2) Re-Fetches model from GDrive
#           2) Saved plot as image
DRAFT_MODE = False  # True = running in background

# To resume the training set the following, remember to set also the correct
epochs= 100 #@param {type:"number"} #100epochs require ~3h for training

## TODO: CHANGE also the reloading code
old = 100 #@param {type:"number"}
#olds = [100, 100, 100, 100, 100]

# 2) Fetching model id of GDrive
MODEL_FILE_GDRIVE_ID = "1wbiFUQIqWlIXkxwHBsX95JuRNmTGgDIX"
BAGGING_MODELS_FILES_GDRIVE_IDS = [
    "1ldeoHu9UxsLI1vidkdpgkqWNp_SxZD0O",#m0
    "1wu15ngwwkJAKlaWYAXz1ZQD3BNM7MF-K", #m1
    "1RjsW44kwyQeBD9KzoVYCX3m_MxnM9Ynv", #m2
    "1U6QqkxZXyTd8kS81ThzdaYRQ53dB7slP", #m3
    "1NZVkPZfOxcshnPU6phyhQMgXkGTocCac", #m4
]
best_val_loss = 10#1.6286746263504028 #@param {type:"number"}
#best_val_loss_epoc = 662 #@param {type:"number"}

PRINT = False

# if you want to save models every 100 epochs
SNAPSHOT = True #@param ["True", "False"] {type:"raw"}
# if you want to save the model on last epoch KEEP IN MIND: vgnet_v3_adam_1e
RESAVE_LAST_SNAPSHOT = True #@param ["True", "False"] {type:"raw"}

#params = [{"MODEL":MODEL},{"OPTIMIZER":OPTIMIZER}]#,LEARNING_RATE, epochs,
#params = [MODEL,OPTIMIZER,LEARNING_RATE, epochs, old, best_val_loss, best_v
print("STARTING WITH params")
#for param in params:
#    param_name = [k for k, v in list(locals().items()) if v == param][0]
#    print(f"  {param_name}:{(16-len(param_name))*' '}{param}")

print(f"  MODEL:{(16-len('MODEL'))*' '}{MODEL}")
print(f"  OPTIMIZER:{(16-len('OPTIMIZER'))*' '}{OPTIMIZER}")
print(f"  LEARNING_RATE:{(16-len('LEARNING_RATE'))*' '}{LEARNING_RATE}")
print(f"  TRAIN:{(16-len('TRAIN'))*' '}{TRAIN}")
```