

# Introduction to Deep Learning (|2DL)

Exercise 10: Semantic Segmentation

### Today's Outline

- Exercise 09: Example Solutions
- Exercise 10: Semantic Segmentation
  - Task & Loss Function
  - Architecture and Upsampling







## **Exercise 9: Solutions**

#### Facial Keypoints

(1, 96, 96) grayscale image

Score: 1/(2\*MSE)

### Threshold: Score of 100 (↔ MSE < 0.005)



#### Leaderboard

#	User	Score	
1	u1180	1584.89	
2	u1345	1361.77	MSE 0.00032
3	u1605	1246.78	
4	u0497	1180.40	
5	u0225	1157.30	
6	u0318	1153.99	
7	u0798	1132.49	
8	u0088	1093.72	
9	u1479	1093.33	
10	u0832	1002.68	
11	u0462	972.42	
12	u0472	924.34	

#### Leaderboard (earlier semester)

#### Leaderboard: Submission 9

Rank	User	Score	Pass
#1	s0672	942.66	1
#2	s0463	940.88	<ul> <li>MSE 0.00053</li> </ul>
#3	s0770	792.80	<ul> <li></li> </ul>
#4	s0303	722.08	1
#5	s0587	689.02	×
#6	s0747	656.89	8°
#7	s0555	654.95	<ul> <li></li> </ul>
#8	s0400	615.63	×
#9	s0322	607.35	*
#10	s0288	602.19	

#### Case Study: Model

```
self.model = nn.Sequential(
    nn.Conv2d(1, 32, (3, 3), stride=1, padding=2),
   # nn.BatchNorm2d(32).
    # nn.Dropout2d(0.2),
   nn.PReLU(),
    nn.MaxPool2d(3),
    nn.Conv2d(32, 64, (3, 3), stride=1, padding=2),
   # nn.BatchNorm2d(64),
    # nn.Dropout2d(),
   nn.PReLU(),
    nn.MaxPool2d(3, stride=2),
    nn.Conv2d(64, 64, (3, 3), stride=1, padding=1),
   # nn.BatchNorm2d(64),
    # nn.Dropout2d(0.3),
   nn.PReLU(),
    nn.MaxPool2d(2, stride=2),
    nn.Conv2d(64, 128, (2, 2), stride=1, padding=1)
   # nn.BatchNorm2d(128),
    # nn.Dropout2d(0.3),
   nn.PReLU(),
```

#### Classic ConvNet architecture:

- Feature extraction
- Classification



#### Case Study: Model Summary

#!pip install torchsummary
import torchsummary

torchsummary.summary(model, (1, 96, 96))

Param #	Output Shape	Layer (type)
320	[-1, 32, 98, 98]	Conv2d-1
1	[-1, 32, 98, 98]	PReLU-2
Θ	[-1, 32, 32, 32]	MaxPool2d-3
18,496	[-1, 64, 34, 34]	Conv2d-4
1	[-1, 64, 34, 34]	PReLU-5
Θ	[-1, 64, 16, 16]	MaxPool2d-6
36,928	[-1, 64, 16, 16]	Conv2d-7
1	[-1, 64, 16, 16]	PReLU-8
Θ	[-1, 64, 8, 8]	MaxPool2d-9
32,896	[-1, 128, 9, 9]	Conv2d-10
1	[-1, 128, 9, 9]	PReLU-11
Θ	[-1, 10368]	Flatten-12
2,654,464	[-1, 256]	Linear-13
0	[-1, 256]	Dropout-14
1	[-1, 256]	PReLU-15
7,710	[-1, 30]	Linear-16

Total params: 2,750,819
Trainable params: 2,750,819
Non-trainable params: 0
······
Input size (MB): 0.04
Forward/backward pass size (MB): 6.72
Params size (MB): 10.49
Estimated Total Size (MB): 17.25

(9x9x128 = 10368)

```
Flatten(),
nn.Linear(10368, 256),
# nn.BatchNorm1d(256),
nn.Dropout(0.1),
nn.PReLU(),
nn.Linear(256, 30),
```

### Case Study: Smaller Linear Layer?

- **1.** Convolutional layer to reduce size to 1x1
  - Here: 9x9 kernel, 128 filters, no padding  $=> 1 \times 1 \times 128 = 128$
- 2. Global Average Pooling (GAP)
  - Here: 9x9 kernel => 128
    - Disadvantage: lose spatial relations
- 3. Flatten

I2DI : Prof. Dai

- Solutions: first use 1x1 convolutions



Extration



FC

Pooling

#### Case Study: With 1x1 Conv

Layer (type)	Output Shape	Param #
Conv2d-1	[-1, 32, 98, 98]	320
PReLU-2	[-1, 32, 98, 98]	1
MaxPool2d-3	[-1, 32, 32, 32]	Θ
Conv2d-4	[-1, 64, 34, 34]	18,496
PReLU-5	[-1, 64, 34, 34]	1
MaxPool2d-6	[-1, 64, 16, 16]	Θ
Conv2d-7	[-1, 64, 16, 16]	36,928
PReLU-8	[-1, 64, 16, 16]	1
MaxPool2d-9	[-1, 64, 8, 8]	Θ
Conv2d-10	[-1, 128, 9, 9]	32,896
PReLU-11	[-1, 128, 9, 9]	1
Conv2d-12	[-1, 16, 9, 9]	2,064
Flatten-13	[-1, 1296]	Θ
Linear-14	[-1, 256]	332,032
Dropout-15	[-1, 256]	Θ
PReLU-16	[-1, 256]	1
Linear-17	[-1, 30]	7,710

# nn.Conv2d(128, 16, (1, 1), stride=1, padding=0),

Total params: 430,451
Trainable params: 430,451
Non-trainable params: 0
Input size (MB): 0.04
Forward/backward pass size (MB): 6.66
Params size (MB): 1.64
Estimated Total Size (MB): 8.34

Next steps:

Make deeper and use residual connection to make it train

# After adding 1x1 layers

#### Case Study: Hyperparameters

```
hparams = {
    "lr": 0.0001,
    "batch_size": 512,
    # TODO: if you have any model arguments/hparams, define them here
}
```

- Default learning rate
- Experiment with batch normalization / Dropout
- Forms of ReLU activations (PReLu, ELU)
- Appropriate weight initialization



## Exercise 10 Semantic Segmentation

#### **Semantic Segmentation**



#### Input: (3xWxH) RGB image

#### Output:

(23xWxH) segmentation map with scores for every class in every pixel

I2DL: Prof. Dai

#### Semantic Segmentation Labels

object class	R	G	В	Colour
void	0	0	0	
building	128	0	0	
grass	0	128	0	
tree	128	128	0	
cow	0	0	128	
horse	128	0	128	
sheep	0	128	128	
sky	128	128	128	
mountain	64	0	0	

"void" for unlabelled pixels







#### **Metrics: Loss Function**

• Averaged per pixel cross-entropy loss



ignore\_index (*int*, optional) – Specifies a target value that is ignored and does not contribute to the input gradient. When size\_average is True, the loss is averaged over non-ignored targets.

#### **Metrics: Accuracy**

• Only consider pixels which are not "void"

```
def evaluate_model(model):
    test_scores = []
    model.eval()
    for inputs, targets in test_loader:
        inputs, targets = inputs.to(device), targets.to(device)
        outputs = model.forward(inputs)
        , preds = torch.max(outputs, 1)
        targets_mask = targets >= 0
        test_scores.append(np.mean((preds == targets)[targets_mask].data.cpu().numpy()))
```

```
return np.mean(test_scores)
print("Test accuracy: {:.3f}".format(evaluate_model(dummy_model)))
```



## **Model Architecture**

### Semantic Segmentation Task

- Input shape: (N, num\_channels, H, W)
   Output shape: (N, num\_classed, H, W)
- We want to:
  - Maintain dimensionality (H, W)
  - Get features at different spatial resolutions



#### **Naive Solution**

- Keep dimensionality constant throughout the network
- Use increasing filter sizes



- Problem:
  - Increased memory consumption
    - Filter size would be the same e.g., 128 filters a (64x3x3) -> 73k params
    - But we have to save inputs and outputs for every layer e.g., 128 filters a (64xWxH) -> millions of params!

#### Excursion: Receptive Field (RF)

- Region in input space a feature in a CNN is looking at
- E.g., after 2 (5x5) convolutions with stride 1 we have a receptive field of 9x9 (RF after first conv: 5 RF after second conv: 5+4)





### **Coming from Classification**

- Use strided convolutions and pooling to increase the receptive field
- Upsample result to input resolution

#### convolution



#### **Better Solution**

- Slowly reduce size -> slowly increase size
  - Pooling -> Upsampling
  - Strided convolution -> Transposed convolution
- Combine with normal convolutions, bn, dropout, etc.



#### **Transposed Convolutions**

• Upsampling with learnable parameters



out = (in - 1) \* stride - 2 \* pad + kernel

(Transpose computation not relevant for the exam, more info here: <a href="https://github.com/vdumoulin/conv\_arithmetic">https://github.com/vdumoulin/conv\_arithmetic</a>)

#### Are transpose convolutions superior?

- Short answer: no, not always
- Long answer: possible checkerboard artifacts for general image generation, see <u>https://distill.pub/2016/deconv-checkerboard/</u>



- My personal go-to:
  - Regular upsampling, followed by a convolution layer

#### How to compete/get results quickly?

• Transfer Learning!



- Possible solutions
  - "The Oldschool"
    - Take pretrained Encoder, set up decoder and only train decoder
    - Encoder candidates: AlexNet, MobileNets
  - "The Lazy"
    - Take a fully pretrained network and adjust outputs

### Summary

- Monday 14.07.25: Lecture 11
  - Guest Lecture
- Wednesday 16.07.25: Exercise 10 Submission
  - Semantic Segmentation: 16.07.25 23:59:59
- Thursday 17.07.25: Tutorial Session 11
- Monday 21.07.25: Lecture 12
  - Recurrent Neural Networks

# Good luck & see you next week 23