

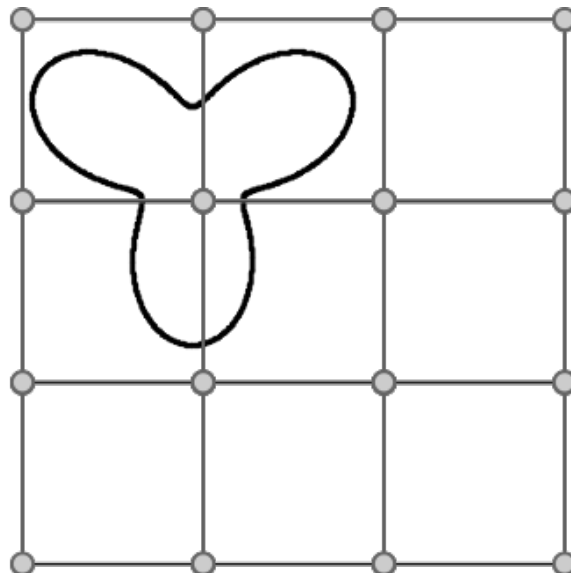
COMPUTER VISION

EXERCISE 5 – NEURAL IMPLICIT MODELS AND RECOGNITION

1 Pen and Paper

1.1 2D Occupancy Network

- a) Derive the decision boundary for a linear layer $L_{\mathbf{W},b}$ with weight $\mathbf{W} = (w_1, w_2)$ and bias b , i.e. where $L_{\mathbf{W},b}(\mathbf{x}) = 0$
- b) Draw the decision boundary for
 - i) $\mathbf{W} = (-5, 1)$, $b = 2$
 - ii) $\mathbf{W} = (2, 3)$, $b = -2$
- c) Now consider a 2-layer fully-connected ReLU-network with weights $\mathbf{W}_1 = \begin{pmatrix} -4 & 1 \\ 2 & 3 \end{pmatrix}$, $\mathbf{b}_1 = \begin{pmatrix} 3 \\ -2 \end{pmatrix}$, and $\mathbf{W}_2 = (2, 3)$, $b_2 = -2$. Let $\sigma(\mathbf{x}) = \begin{pmatrix} \max(0, x_1) \\ \max(0, x_2) \end{pmatrix}$ denote the ReLU activation function. Derive and draw the decision boundary for this network. You can use an online tool, e.g., <https://www.geogebra.org/m/Adc44ZZq> to plot the decision boundary.
- d) So far, we have learned how to extract the decision boundary from a trained network. Next, we want to extract a mesh from this implicitly defined boundary. The figure below shows the decision boundary for some classifier. Apply the MISE algorithm manually by coloring occupied points in red and free points in blue and refining the grid where needed. See Lecture 9, Slide 13 for an example. Refine the grid for 3 steps. Lastly, manually execute the marching squares algorithm, the 2D variant of the marching cubes algorithm, and draw the extracted mesh into the last refined grid. Assume that the classifier returns either 0 (free space) or 1 (occupied).



1.2 Differentiable Volumetric Rendering

Let the occupancy network f_{θ_f} in DVR be parametrized by a single linear layer, i.e. $f_{\theta_f} = L_{\mathbf{W}_f, b_f}$. Further, let the texture field t_{θ_t} be defined by a linear layer with a subsequent sigmoid function σ , i.e. $t_{\theta_t} = \sigma(L_{\mathbf{W}_t, b_t})$. Similarly to the previous exercise, we choose the weights: $\mathbf{W}_f = (-4, 1)$, $b_f = 3$, and $\mathbf{W}_t = (2, 3)$, $b_t = -2$. Note that in this case the texture field only predicts a single scalar that corresponds to the intensity of the pixel to simplify computations.

- a) We shoot a ray from the origin in an angle of 45° , i.e. $\mathbf{r} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} + \hat{d} \begin{pmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{pmatrix}$. Estimate the surface intersection point $\hat{\mathbf{p}}$ by performing a single update step with the secant method starting in the interval $[\hat{d}_1 = 1, \hat{d}_2 = 2]$.
- b) Compute the predicted texture field, i.e. intensity, at $\hat{\mathbf{p}}$.
- c) As loss \mathcal{L} we choose the L2-distance between the predicted intensity \hat{I} and the ground truth value I . Let us assume that $I = 0.6$. The gradient wrt. all parameters $\theta = \theta_f \cup \theta_t$ is

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \theta} &= \frac{\partial \mathcal{L}}{\partial \hat{I}} \frac{\partial \hat{I}}{\partial \theta} = \frac{\partial \mathcal{L}}{\partial \hat{I}} \frac{dt_{\theta_t}(\hat{\mathbf{p}})}{d\theta} \\ &= \frac{\partial \mathcal{L}}{\partial \hat{I}} \left[\frac{\partial t_{\theta_t}(\hat{\mathbf{p}})}{\partial \theta} + \frac{\partial t_{\theta_t}(\hat{\mathbf{p}})}{\partial \hat{\mathbf{p}}} \frac{\partial \hat{\mathbf{p}}}{\partial \theta} \right] \end{aligned}$$

Therefore the gradients wrt. the parameters f_{θ_f} , i.e. \mathbf{W}_f and b_f are given by:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_f} = \frac{\partial \mathcal{L}}{\partial \hat{I}} \left[\frac{\partial t_{\theta_t}(\hat{\mathbf{p}})}{\partial \mathbf{W}_f} + \frac{\partial t_{\theta_t}(\hat{\mathbf{p}})}{\partial \hat{\mathbf{p}}} \frac{\partial \hat{\mathbf{p}}}{\partial \mathbf{W}_f} \right] \quad (1)$$

$$\frac{\partial \mathcal{L}}{\partial b_f} = \frac{\partial \mathcal{L}}{\partial \hat{I}} \left[\frac{\partial t_{\theta_t}(\hat{\mathbf{p}})}{\partial b_f} + \frac{\partial t_{\theta_t}(\hat{\mathbf{p}})}{\partial \hat{\mathbf{p}}} \frac{\partial \hat{\mathbf{p}}}{\partial b_f} \right] \quad (2)$$

Start by computing the following components of Eq. (1),(2) $\frac{\partial \mathcal{L}}{\partial \hat{I}}$, $\frac{\partial t_{\theta_t}(\hat{\mathbf{p}})}{\partial \mathbf{W}_f}$, $\frac{\partial t_{\theta_t}(\hat{\mathbf{p}})}{\partial b_f}$, $\frac{\partial t_{\theta_t}(\hat{\mathbf{p}})}{\partial \hat{\mathbf{p}}}$.

- d) To compute the gradients for \mathbf{W}_f and b_f the only components missing for Eq. (1),(2) are $\frac{\partial \hat{\mathbf{p}}}{\partial \mathbf{W}_f}$ and $\frac{\partial \hat{\mathbf{p}}}{\partial b_f}$, respectively. However, we cannot directly compute them because $\hat{\mathbf{p}}$ is defined only implicitly. Use implicit differentiation on $f_{\theta_f} = \tau$, where $\tau = 0$, to obtain $\frac{\partial \hat{\mathbf{p}}}{\partial \mathbf{W}_f}$, $\frac{\partial \hat{\mathbf{p}}}{\partial b_f}$. With these results, compute the gradients for \mathbf{W}_f and b_f .

Hint: Write $\hat{\mathbf{p}}$ as the ray from the origin $\hat{\mathbf{p}} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} + \hat{d} \begin{pmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{pmatrix}$ to link it to the depth prediction \hat{d} .

Since the depth is predicted by f_{θ_f} , this will allow you to formulate $\frac{\partial \hat{\mathbf{p}}}{\partial \mathbf{W}_f}$, $\frac{\partial \hat{\mathbf{p}}}{\partial b_f}$ in terms of $\frac{\partial f_{\theta_f}(\hat{\mathbf{p}})}{\partial \hat{\mathbf{p}}}$ and $\frac{\partial f_{\theta_f}(\hat{\mathbf{p}})}{\partial \mathbf{W}_f}$ and $\frac{\partial f_{\theta_f}(\hat{\mathbf{p}})}{\partial b_f}$, respectively. For computing the result this means you also need to compute $\frac{\partial f_{\theta_f}(\hat{\mathbf{p}})}{\partial \hat{\mathbf{p}}}$, $\frac{\partial f_{\theta_f}(\hat{\mathbf{p}})}{\partial \mathbf{W}_f}$ and $\frac{\partial f_{\theta_f}(\hat{\mathbf{p}})}{\partial b_f}$.

1.3 Neural Network Layers

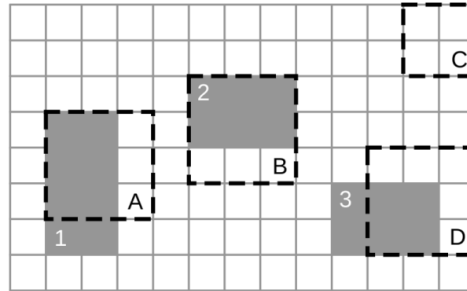
- a) Calculate the number of parameters $\#P$ and the receptive field size S_r for the following layers
 - i) Linear layer with input dimension $D_{in} = 16^2$ and output dimension $D_{out} = 256$ and bias and a flattened grayscale image as input.
 - ii) Convolution with filter size $K = 3$, $C_{in} = 32$ input and $C_{out} = 64$ output channels, padding $P = 1$, stride $S = 1$ and bias
 - iii) Two convolutions, both with filter size $K = 3$, $C_{in} = 32$ input and $C_{out} = 32$ output channels, padding $P = 1$, Dilation $D = 2$, stride $S = 1$ and bias

- iv) Two convolutions, both with filter size $K = 3$, $C_{in} = 32$ input and $C_{out} = 32$ output channels, padding $P = 1$, Dilation $D = 1$, stride $S = 2$ and bias
- v) Average Pooling with size 2
- b) Complete the following table for softmax classification with Cross Entropy Loss. Each row is a different sample of a minibatch and the values of the ground truth class are marked in **red**. Which samples are classified correctly? Which samples are classified incorrectly? Which samples contribute most strongly to the parameter update of the model?

| Predicted scores \mathbf{s} | $\text{softmax}(\mathbf{s})$ | CE loss |
|-------------------------------|------------------------------|---------|
| $(-1, +3, +0, +1)^T$ | | |
| $(+2, +3, -1, -1)^T$ | | |
| $(+0, +0, +3, +4)^T$ | | |
| $(+1, +0, +3, +4)^T$ | | |

1.4 Detection Metrics

Consider the following image where the three gray shaded boxes (1,2,3) correspond to the true objects and the four dashed boxes (A,B,C,D) correspond to detections of an object detector:



- a) Calculate the IoU (Intersection-over-Union) metric for the pairs (A, 1), (B, 2) and (D, 3). Assuming a detection threshold of $\tau = 0.5$, which of the three objects have been correctly detected?
- b) Calculate the number of true positives (TP), false positives (FP) and false negatives (FN) assuming the same detection threshold as in the previous question ($\tau = 0.5$).
- c) Calculate precision and recall at a detection threshold of $\tau = 0.6$.

2 Coding Exercises

The following two exercises will be coding exercises. They are contained in the jupyter notebooks `code/implicit.ipynb` and `code/recognition.ipynb`. The notebooks are self-contained but you can also use this document as guidance. If you are stuck, you can find Hints in the notebook itself which are written upside-down.

2.1 Occupancy Network

In this exercise, you will train an occupancy network to represent a single 3D object. More specifically, you will start from the point cloud of a chair together with its occupancy values, see Fig. 1. Then, you will train an occupancy network to encode the shape and reconstruct the mesh from the trained network. The code is provided in `code/implicit.ipynb`. Make sure to copy the data for this task `code/data/points.npz` into your working directory to a folder called `data` if you use Google Colab.

- a) Prepare your data by implementing the function `get_train_val_split`.

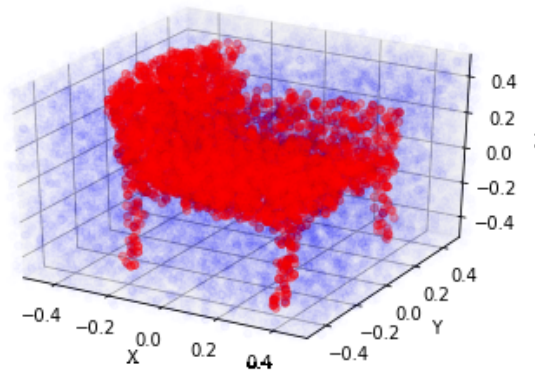


Figure 1: **3D Point Cloud of a Chair.** Red points belong to the object while blue points represent free space.

- b) Complete the class `OccNet` such that you obtain a network with 4 hidden layers of hidden dimension 64. The network takes a batch of 3D locations as input and predicts a single occupancy value for each point.
- c) Complete the training loop, i.e. the function `train_model`. Do the curves for training and validation loss look good? Or is there under / overfitting?
- d) Load the trained model and predict the occupancy values for all points on a 3D grid. You can use the `test_loader` to load the grid points in batches.

2.2 Person Detector

In this exercise, you will develop your own person detector. You will extract your training data using images of pedestrians with bounding box labels, see Fig. 2. Then, you train two simple classifiers on the feature representation of the training data. Lastly, you use the trained classifiers to detect pedestrians in new images using a sliding window approach. The code is provided in `code/recognition.ipynb`. Extract the dataset in `code/data/PennFudanPed.zip` using `unzip code/data/PennFudanPed.zip -d code/data/`.

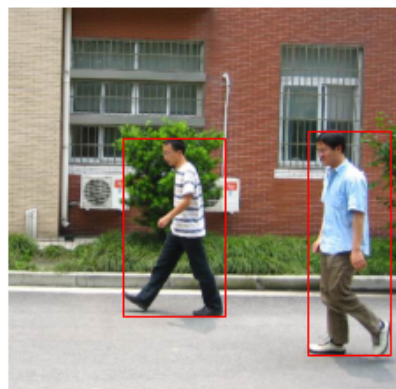


Figure 2: **Person Detector.** Example image of pedestrians with bounding box labels.

- a) To mine negative training samples, implement the function `get_random_box`. The function should return the (left, upper, right, lower) pixel coordinates of a box of a given size at a random location in the image.

- b) To further improve the accuracy our classifier we can add extra difficult negative samples to our training data. We construct hard negative examples by adding a small offset to the ground truth boxes. Thereby, the negative examples are close to positive detections. Implement the function `add_offset_to_box`. This function takes a box and an image size as input and returns a box that is slightly translated in the image.
- c) Implement the function `img_to_hog_patches` that extracts image patches of a given size and converts these patches to HOG features. Of course, you should use the provided helper functions, e.g. to get the HOG features.