

# Object Recognition and Image Understanding

## Exercise Sheet 6

Shivali Dubey, Sascha Stelling

June 8, 2018

### Question 1

- **Team:**

Shivali:

- Implementation
- Neural network setup

Sascha:

- Implementation
- Preparation of the dataset
- Running training phase

- **Problem Definition:** *Object detection and image classification* has various applications, for instance, in self-driving cars to detect and classify pedestrians, motorcycles, trees, bicycles etc; classification of features on the Earth such as roads, rivers, agricultural fields etc using satellite images. With the advancements in deep learning, every year, new algorithms/models keep on outperforming the previous ones, to achieve the best possible accuracies for image classification. One of the most popular dataset used is the ImageNet dataset. In our project we propose to implement the deep learning algorithms based on a few selected studies [1,2,3] with the aim to attain the best possible classification accuracy.
- **Dataset:** Tiny ImageNet<sup>1</sup>. Tiny Imagenet has 200 classes. Each class has 500 training images, 50 validation images, and 50 test images.
- **Approach:** With such a large dataset, one of the main challenges of classification is diversity of the images. Our model/algorithm must be able to handle fine-grained and specific classes even when they are hard to distinguish. In other words, we need to maximize inter-class variability, while

---

<sup>1</sup><https://tiny-imagenet.herokuapp.com/>

minimize intra-class variability. At the same time, attaining the best possible classification accuracy is always a challenge for any given algorithm. The predictions go wrong when you have too many false positives and false negatives.

Image Classification:  
Architecture:

- Convolution: The main purpose of using multiple convolution layers is feature extraction. We use Scale-Invariant Feature Transform (SIFT)<sup>2</sup> descriptors which computes the Difference of Gaussians (DoG)<sup>3</sup>. DoG is used to detect blobs by subtracting two blurred images from another with different Gaussian kernels. The maxima and minima of this operation are taken by SIFT as key feature locations for the next neurons. To classify feature more accurately we make use of densely sampled SIFT, Extended Opponent SIFT and RGB-SIFT detector as described in [3] in three different convolution layers. We use Parametric ReLU (PReLU)<sup>4</sup> for activation of neuron and Adam<sup>5</sup> for optimization. PReLU is being used instead of ReLU because it improves model fitting with nearly zero computational cost and little overfitting risk. Also, PReLU (and also ReLU) brings non-linearity into the system which allows learning complex functions. The weight initialization can be performed using Xavier's initialization<sup>6</sup>. Furthermore, the weight optimization is also controlled by PReLU.
- Downsampling: Though the activation function passes only relevant pixels to the next layer, the array could still be big. To reduce the size of the array, we downsample it using an algorithm called max pooling.
- Fully-connected Neural Network: We construct the last layer as fully connected neural network with hidden layer and logistic regression, and set all feature maps produced from previous layers as inputs. Softmax logistic regression can be used to represent categorical distribution i.e. probability distribution over different outcomes.
- Backpropagation: We use gradient descent method for optimization algorithm which is thus used for learning and training.

---

<sup>2</sup><https://pdfs.semanticscholar.org/presentation/e903/196678c93315f2bf6f0235b3bab59c157b04.pdf>

<sup>3</sup><http://micro.magnet.fsu.edu/primer/java/digitalimaging/processing/diffgaussians/index.html>

<sup>4</sup>[https://www.cv-foundation.org/openaccess/content\\_iccv\\_2015/papers/He\\_Delving\\_Deep\\_into\\_ICCV\\_2015\\_paper.pdf?spm=5176.100239.blogcont55892.28.pm8zm1&file=He\\_Delving\\_Deep\\_into\\_ICCV\\_2015\\_paper.pdf](https://www.cv-foundation.org/openaccess/content_iccv_2015/papers/He_Delving_Deep_into_ICCV_2015_paper.pdf?spm=5176.100239.blogcont55892.28.pm8zm1&file=He_Delving_Deep_into_ICCV_2015_paper.pdf)

<sup>5</sup><https://arxiv.org/pdf/1412.6980.pdf>

<sup>6</sup><http://proceedings.mlr.press/v9/glorot10a/glorot10a.pdf>

- **Evaluation & Expected Results:** Reducing Overfitting: As studied in previous research works<sup>7</sup>, we perform data augmentation and dropout to reduce overfitting. Data augmentation refers to artificially enlarging image size using augmentation. Dropout refers to dropping out the output of each hidden neuron with probability 0.5, such that the respective neuron can't participate in backpropagation.  
Training: The training is carried out by optimising the multinomial logistic regression objective function using mini-batch gradient-descent (based on back-propagation)(Lecunn, imagenet class).  
We try to achieve at least 70% accuracy with our algorithm.
- **Hardware:** We plan to use our own PCs for this task having the following specifications: Intel Core i5-3350P @ 3.1GHz, 12GB DDR3, Nvidia GTX 970, 128GB SSD
- **Excluded Presentation Date:** None

---

<sup>7</sup><http://cs229.stanford.edu/proj2014/Xiaodong%20Zhou,%20Supervised%20DeepLearning%20For%20MultiClass%20Image%20Classification.pdf>