

## Computer Vision Assignment - 2

### **Q1 Compare and contrast the computation of sliding window Detection and convolution neural networks for object detection.**

**Sol.-**

Sliding window detection and convolutional neural networks (CNNs) are two distinct approaches used for object detection, each with its own set of advantages and limitations. Let's compare and contrast them:

Approach:

- i. **Sliding Window Detection:** This approach involves systematically scanning an image with a fixed-size window, often referred to as the "sliding window," and classifying each window region using a pre-trained classifier. The window slides across the entire image, and for each position, the classifier determines whether an object is present or not.
- ii. **Convolutional Neural Networks:** CNNs are a type of deep learning model specifically designed to automatically and hierarchically learn features from raw data. For object detection, CNNs process the entire image at once through a series of convolutional layers, which extract features at different spatial levels, followed by fully connected layers for classification and localization.

Feature Extraction:

- i. **Sliding Window Detection:** Features are typically handcrafted and provided as input to a classifier. Common choices include Histogram of Oriented Gradients (HOG), Scale-Invariant Feature Transform (SIFT), or Local Binary Patterns (LBP). These features capture basic visual patterns in the image.
- ii. **Convolutional Neural Networks:** CNNs automatically learn features from the data during the training process. The network learns to extract relevant features directly from raw pixels, which can capture intricate details and hierarchical representations of objects.

Localization:

- i. **Sliding Window Detection:** Object localization involves determining the precise bounding box of the detected object within the image. This process is separate from the detection stage and typically requires additional post-processing steps to refine the localization accuracy.
- ii. **Convolutional Neural Networks:** CNNs often include mechanisms for both detection and localization within the same architecture. Techniques like Region Proposal Networks (RPNs) or bounding box regression layers are commonly employed to predict both the presence of objects and their bounding boxes simultaneously.

Training Data and Complexity:

- i. **Sliding Window Detection:** Training a sliding window detector requires a large dataset of labeled examples and extensive tuning of parameters like window size and step size. This approach can become computationally expensive, especially for high-resolution images or large datasets.

- ii. Convolutional Neural Networks: CNNs require labeled training data for supervised learning, but they can automatically learn relevant features and spatial hierarchies from the data, reducing the need for manual feature engineering. While training CNNs can be computationally intensive, once trained, they can efficiently process images of various sizes without significant modifications.

Performance and Flexibility:

- i. Sliding Window Detection: Sliding window detection can be effective for detecting objects of various scales and aspect ratios. However, it may struggle with computational efficiency, especially when dealing with large images or datasets.
- ii. Convolutional Neural Networks: CNNs have demonstrated superior performance in object detection tasks, particularly in terms of accuracy and speed. They can adapt to various object scales and orientations through techniques like multi-scale feature extraction and anchor boxes.

In summary, while sliding window detection offers a straightforward approach to object detection, CNNs provide a more robust and flexible solution by automatically learning relevant features and spatial hierarchies from the data. However, CNNs often require more computational resources for training and inference.

## **Q2 Analysis the differences and similarities of features and classification of methods: HOG, SVM, Adaboosts, CNNs, VGGNet, GoogleNet, Resnet.**

**Sol.**

The differences and similarities in features and classification methods of the following techniques: Histogram of Oriented Gradients (HOG), Support Vector Machines (SVM), Adaboost, Convolutional Neural Networks (CNNs), VGGNet, GoogLeNet, and ResNet.

Histogram of Oriented Gradients (HOG):

Features: HOG extracts local gradient orientation information from image patches. It divides the image into small cells, computes gradient orientations within each cell, and then constructs histograms of gradient orientations across cells.

Classification: HOG features are often used with classifiers like SVM or Adaboost to classify objects in images based on the extracted features.

Support Vector Machines (SVM):

Classification: SVM is a supervised learning algorithm used for classification tasks. It learns a decision boundary that best separates different classes in the feature space. SVMs are commonly used in conjunction with handcrafted features like HOG for object detection tasks.

Adaboost:

Classification: Adaboost (Adaptive Boosting) is an ensemble learning method that combines multiple weak classifiers to create a strong classifier. It sequentially trains a series of weak classifiers, with each subsequent classifier focusing more on the training instances that the previous ones misclassified.

Similarity: Adaboost is often used in conjunction with weak classifiers, including classifiers based on handcrafted features like HOG, to improve classification performance.

#### Convolutional Neural Networks (CNNs):

Features: CNNs automatically learn hierarchical representations of features from raw image data. They consist of multiple layers, including convolutional layers, pooling layers, and fully connected layers. The convolutional layers extract features from the input images by convolving learnable filters over the input.

Classification: CNNs perform end-to-end learning, where both feature extraction and classification are integrated into the network. They have shown remarkable success in various computer vision tasks, including object detection, by learning features directly from the data.

#### VGGNet:

Features: VGGNet is a deep CNN architecture known for its simplicity and uniformity. It consists of multiple convolutional layers, followed by max-pooling layers, and finally fully connected layers. VGGNet primarily focuses on learning deep hierarchical features from images.

Classification: VGGNet can be trained for image classification tasks and can also be adapted for object detection by combining it with region proposal techniques and classification layers.

#### GoogLeNet (Inception):

Features: GoogLeNet introduced the Inception module, which consists of multiple parallel convolutional layers of different kernel sizes and a pooling layer. This architecture aims to capture both local and global features effectively.

Classification: GoogLeNet is primarily used for image classification tasks. Its deep architecture allows it to learn rich representations of images, making it suitable for object detection when combined with appropriate techniques.

#### ResNet:

Features: ResNet introduced residual connections, which facilitate training deeper neural networks by alleviating the vanishing gradient problem. These connections allow information from earlier layers to bypass several layers, enabling the network to learn more robust features.

Classification: ResNet is widely used for image classification tasks due to its ability to train very deep networks effectively. It can also be adapted for object detection tasks by integrating it with region proposal networks and classification layers.

#### Differences and Similarities:

Features: HOG, SVM, and Adaboost rely on handcrafted features extracted from images, whereas CNNs, VGGNet, GoogLeNet, and ResNet learn features directly from raw data.

**Classification:** HOG, SVM, and Adaboost use traditional machine learning classifiers, while CNNs and deep learning architectures like VGGNet, GoogLeNet, and ResNet perform end-to-end learning for classification.

**Complexity:** Deep learning architectures like CNNs, VGGNet, GoogLeNet, and ResNet are generally more complex than traditional methods like HOG, SVM, and Adaboost. However, they often achieve superior performance, especially on large-scale datasets.

**Training:** Traditional methods like HOG, SVM, and Adaboost require manual feature engineering and tuning, while deep learning architectures automatically learn features during training.

**Performance:** Deep learning architectures, particularly CNNs, VGGNet, GoogLeNet, and ResNet, have demonstrated state-of-the-art performance on various computer vision tasks, including object detection, due to their ability to learn complex features and hierarchies directly from data.