
BBM418 Introduction to Computer Vision Lab. # Single Object Tracking with Regression Networks

Meryem Ülkü Kara
21727355
Department of Computer Engineering
Hacettepe University
Ankara, Turkey
b21727355@cs.hacettepe.edu.tr

Overview

In this assignment, it is expected to implement a basic single object viewer by training the network on the given videos. We try to predict the ground truth of moving objects. I handled my homework in two ways.

1 Model

In the first, I grouped the frames and ground truth tables in pairs, passed the steps of reshaping on them, and gave these four data to the model. In the model, I first cropped for the first frame and extracted the feature, then I cropped and extracted the feature for the second frame. I used the first of my ground truth tables for these crops. Then I combined these two properties and aimed to find a prediction with it, it should be predicting the second one of my prediction ground truth tables I found. But while the shape of my predictionum is a value like [1,2000], my ground truth table has only 4 values. There was a size mismatch here and no matter how hard I tried, I couldn't find any solution to it.

```
DATASET: dataset = torch.utils.data.TensorDataset
```

```
(frame1tensor, frame2tensor, boundingbox1tensor, boundingbox2tensor)
```

```
MODEL TRAIN:
```

```
for frame1, frame2, boundingbox1, boundingbox2 in loader : optimizer.zero_grad()
```

```
resizedt transform = transforms.Resize((224, 224))
```

```
CROP FRAME1 cropframe1 = cropframe(frame1, boundingbox1) FEATURE EXTRACTION features1 =  
model(cropframe1)
```

```
CROP FRAME2 cropframe2 = cropframe(frame2, boundingbox1) FEATURE EXTRACTION features2 =  
model(cropframe2)
```

```
COMBINING FEATURE VECTORS combinedfeatures =  
torch.cat((features1, features2), dim = 1)
```

```
print(combinedfeatures.shape)
```

```
Predictions predictions = model.fc(combinedfeatures)
```

```
lass calculation and back propagation loss = criterion(predictions, boundingbox2)
```

```

CROP METHOD : def crop_frame(frame, bbox) :
BOUNDING BOX VALUES x, y, width, height = bbox.squeeze().tolist()
FIND CENTER center_x = x + (width/2) center_y = y + (height/2)
ENLARGED BBOX VALUES enlarged_width = width * 2 enlarged_height = height * 2
CROPPED FRAME cropped_frame = frame[:, :, int(center_y - (enlarged_height/2)) :
int(center_y + (enlarged_height/2)), int(center_x - (enlarged_width/2)) : int(center_x +
(enlarged_width/2))] return cropped_frame

```

2 Model2

In the second time, I read the frames and ground truths and grouped them in pairs. Since there was a size mismatch in the model and I couldn't solve it in my first handling, I cropped my frames as desired before entering the model in this model and found the expected cropped frames. Then I got my cropped frame1, cropped frame2 , expected bounding box , expected cropped frame2. In this part, I tried to train the model in two ways, the first was to train the cropped frame1, cropped frame2 , and the expected cropped frame2 data. I aimed to extract and combine features for Frame1 and Frame2, then make predictions with this feature extraction and find the expected frame2, but here again I encountered a size mismatch. I tried many things like removing fc-layer but I couldn't get past this issue. On the other hand, I tried to find the expected bounding box by taking the cropped frame1, cropped frame2 data as input value, but I couldn't get over the size mismatch errors.

CROP METHOD :

```

def crop_frame(frame, bbox) :
GET BBOX VALUES x, y, width, height = bbox
frame_height, frame_width, = frame.shape
CONTROL BOUNDING BOX HEIGHT LARGER OR NOT TO FRAME HEIGHT if height >
frame_height : height = frame_height CONTROLFORWEIGHT if width > frame_width :
width = frame_width
CALCULATE CENTER center_x = x + (width/2) center_y = y + (height/2)
CALCULATE ENLARGED BBOX enlarged_width = width * 2 enlarged_height = height * 2
FIND CROPPED FRAME cropped_frame = frame[int(center_y - (enlarged_height/2)) :
int(center_y + (enlarged_height/2)), int(center_x - (enlarged_width/2)) : int(center_x +
(enlarged_width/2))]
return cropped_frame

```

2.1 Model2-1

DATASET: dataset = CustomDataset(cropped_images1, cropped_images2, cropped_images3)

MODEL TRAIN

for frame1, frame2, frame3 in data_loader : optimizer.zero_grad()

EXTRACT FEATURES WITH CROPPED_FRAME1 features1 =
model(frame1).Featureextraction

EXTRACT FEATURES WITH CROPPED_FRAME2 features2 =
model(frame2).Featureextraction

COMBINING FEATURES VECTORS combined_features =
torch.cat((features1, features2), dim = 1)

```
print(combined_features.shape)print(frame3.shape)
Predictions predictions = combined_featuresprint(predictions)
lass calculation and back propagation loss = criterion(combined_features, frame3)loss.backward()
```

2.2 Model2-2

```
DATASET: dataset = CustomDataset(cropped_images1, cropped_images2, bbox)
MODEL TRAIN feature_extractor = torch.nn.Sequential(*list(model.children())[:-1])
for frames, targets in data_loader :
    frame1 = frames[0].to(device) frame2 = frames[1].to(device) targets = targets.to(device)
    optimizer.zero_grad()
    FEATURE EXTRACT features1 = feature_extractor(frame1) features2 =
    feature_extractor(frame2)
    COMBINE FEATURE VECTORS combined_features = torch.cat((features1.flatten(start_dim =
    1), features2.flatten(start_dim = 1)), dim = 1)
    hidden = torch.relu(fc1(combined_features)) predictions = fc2(hidden)
    PREDICTED GROUND TRUTH BOXES ground_truth_box = frame3
    loss = criterion(predictions, ground_truth_box)loss.backward()optimizer.step()

I tried many methods, but in all of them I got errors during the training phase and could not solve my
mistakes.
```