

Tennis Shot Classification: Using Machine Learning to Classify Types of Tennis Shots in Video Data

Irene Agusti, Utsah Kalra, Sam Rudenberg

May 17, 2024

Abstract

We created machine learning models to predict the type of tennis shot in a given frame of a video. Video data was split into sequential frames, which were manually labeled by the type of shot discernible by the human eye. We created a CNN model, and attempted to create a 3DCNN model, and evaluated them on training and testing data. Our results found that our CNN model reached an accuracy of 17% on the training set, and 19% on the test set.

Intro

In this project, we use a data set containing snapshots from a tennis match to classify the different types of tennis shots being hit in that snapshot. There are many machine learning algorithms that one can use to do this. In our project, we choose to focus on using Convolutional Neural Networks to perform this classification task. While CNN's are computationally expensive, they are well-suited to data that requires identification and extraction of spatially dependent hidden features (slides).

Methods

The data used for the project comes from Github user yastrebksy. This user has a repository for TennisProject, which is described as Tennis analysis using deep learning and machine learning. A subcomponent of this user's project, TrackNet is used to take video data as input, and output frame by frame information with file names, visibility, xy-coordinates, and status. File names are individually labeled frames. Visibility describes how visible a tennis ball is within a given frame; 0 for ball not in image, 1 for ball can be easily identified, 2 for in frame, but not easily identified. xy-coordinates are provided for the location of the ball in the frame (where 0,0 is the bottom left corner). Status has three options: hitting, flying, or bouncing.

We downloaded images and outputs from TrackNet for Games 1 and Games 8. The images we downloaded are pre-cut, frame by frame pieces of a 30-fps video clip. After downloading all the images, we wrote code to rename them, appending the game and clip number, along with the sequence number of the image where 0 is the first frame of the video, and the largest number is the last frame of the video. We added a column to the TrackNet output called shot, with the following classifications: 1 - forehand, 2 - backhand, 3 - forehand slice, 4 - backhand slice, 5 - serve, 6 - forehand volley, 7 - backhand volley, 8 - overhead, 9 - standby. After all data was cleaned and labeled, we moved on to the model. Below is a breakdown of the frames for each clip in game 1 as well as game 8. We renamed the images as seen in the sequenced columns so that all images could be in the same folder, while having a way of being differentiated from the other clips. This was a replacement for the earlier game.clip.frame appendage that we used for our 3DCNN model.

Game 1, Clip:	Sequenced	Total Frames
1	1-207	207
2	208-337	130
3	338-400	63
4	401-508	108
5	509-668	160
6	669-746	78
7	747-1198	452
8	1199-1313	115
9	1314-1422	109
10	1423-1552	130
11	1553-1678	126
12	1679-1931	253

Game 8, Clip:	Sequenced	Total
1	1932-2049	118
2	2050-2144	95
3	2145-2527	383
4	2528-2660	133
5	2661-2739	79
6	2740-2805	66
7	2806-3227	422
8	3228-3546	319
9	3547-3784	238

In the context of our model, our predictors/features are the images and the targets/classification are the shot type labels, which we created ourselves. Our image data size was started at [720, 1280, 3].

We began to create our model. Our first model was a CNN image classification model with five layers: Convolutional, ReLU, Max Pool, and two linear layers. The first layer was the convolutional layer, which consisted of 3 RGB inputs, 16 feature maps, a kernel size of 5, stride of 1, and padding of 2. The next layer was the ReLU layer, which acted as an activation layer. This was followed by a Max Pool layer with a kernel of 2, and stride of 2, which reduced the size of the feature maps by taking the max value. Lastly, we have two linear layers which reduce the inputs from 134400 to 128, to 9 outputs which serve as our classifications. This was presented in class.

When we first ran our training data on the model, we encountered issues with RAM. Our image size was too large, and was causing our computers to run out of memory. To counteract this, we decided to reduce our image size. We used interpolation to reduce our image size to 140x240. By reducing the image this much, there is a risk of losing quality, but it was necessary to make sure our computer's had enough RAM to run the model. After running the model on the training data, we evaluated the model using cost curves and classification accuracy. Our initial accuracy was 20%.

After our class presentation, we discovered that the data set we had downloaded from the Github project had incomplete games, which meant our results were invalid. We had originally used clips from game 9, rather than the clips from game 8 which are seen above, but this meant we had to evaluate our model.

As we re-evaluated our model, it was even more demanding in terms of RAM. We converted our images to grayscale, and simplified our model. We experimented with different batch sizes, learning rates, and framework layers to maximize accuracy while still being small enough for our computer to run. Our network became a series of three linear layers with ReLU activation between them, reducing from 33600 inputs, to 8192, to 512, to 10. We used 5 epochs, a learning rate of 0.01, and a batch size of 100.

While a CNN model allows us to exploit the spatially dependent features in image data, it does not allow us to exploit the sequential nature of data. Here in our project, the data set was created from video data and has an important sequential aspect to it. Using Recurrent Neural

Networks allows us to take into account the sequential aspect of our data set. The mechanism through which rnn's performs sequential learning is the "hidden state." This "hidden state" maintains information about the previous frame, and uses both the frames features itself, combined with the information from the previous frame to predict the next outcome.

More specifically, for our second comparative model using recurrent neural networks, we use a LSTM (Long Short Term Memory) + CNN model. The LSTM+ CNN model handles the sequential learning and relative feature learning aspects separately. CNN handles the spatial feature extraction while the LSTM focuses on temporal relationships. This separation can make it easier to tune to each part of the model.

Another way to take into account the sequential nature of data is to perform a 3D CNN model. A 3D CNN model, treated the sequence of frames as a third layer dimension of the data instead of dealing with it separately as in the LSTM+CNN model. For this model, we have a convolutional layer with 16 feature maps, a kernel size of 3, a stride of 3 and padding of 1; a max pooling layer with kernel 2, and stride of 2. Finally, we have 2 fully connected layers, the first one takes an input of 134,400 and compresses it to output 128 features; the second connected layer reduces it to the final 10 outputs that serve as our classifications.

Results

After 1 epoch, our model could do no further learning due to the simplicity of our model. Our final CNN model made correct predictions on the training set 17% of the time, and 19% of the time on the test set. Due to the complexities and time restrictions, we were unable to run our other models. Below are epochs, epoch images from running the model.

Epoch 1/5:	3%	1/35 [00:03<01:52, 3.32s/batch]Cost: 2.33108543496704
Epoch 1/5:	6%	2/35 [00:06<01:41, 3.07s/batch]Cost: 2.4317221641540527
Epoch 1/5:	9%	3/35 [00:09<01:37, 3.06s/batch]Cost: 2.2811508178710938
Epoch 1/5:	11%	4/35 [00:13<01:47, 3.48s/batch]Cost: 2.3011507987976074
Epoch 1/5:	14%	5/35 [00:16<01:38, 3.28s/batch]Cost: 2.331150770187378
Epoch 1/5:	17%	6/35 [00:19<01:31, 3.15s/batch]Cost: 2.26115083694458
Epoch 1/5:	20%	7/35 [00:22<01:33, 3.34s/batch]Cost: 2.271150827407837
Epoch 1/5:	23%	8/35 [00:26<01:33, 3.47s/batch]Cost: 2.2911508083343506
Epoch 1/5:	26%	9/35 [00:30<01:30, 3.49s/batch]Cost: 2.2811508178710938
Epoch 1/5:	29%	10/35 [00:33<01:21, 3.28s/batch]Cost: 2.2811508178710938
Epoch 1/5:	31%	11/35 [00:37<01:20, 3.68s/batch]Cost: 2.321150779724121
Epoch 1/5:	34%	12/35 [00:41<01:23, 3.64s/batch]Cost: 2.3511507511138916
Epoch 1/5:	37%	13/35 [00:44<01:21, 3.70s/batch]Cost: 2.2511508464813232
Epoch 1/5:	40%	14/35 [00:47<01:12, 3.43s/batch]Cost: 2.3111507892608643
Epoch 1/5:	43%	15/35 [00:50<01:04, 3.24s/batch]Cost: 2.2811508178710938
Epoch 1/5:	46%	16/35 [00:53<00:59, 3.11s/batch]Cost: 2.3511507511138916
Epoch 1/5:	49%	17/35 [00:57<00:59, 3.33s/batch]Cost: 2.2311508655548096
Epoch 1/5:	51%	18/35 [01:00<00:57, 3.40s/batch]Cost: 2.3411507606506348
Epoch 1/5:	54%	19/35 [01:03<00:52, 3.25s/batch]Cost: 2.331150770187378
Epoch 1/5:	57%	20/35 [01:06<00:46, 3.13s/batch]Cost: 2.331150770187378
Epoch 1/5:	60%	21/35 [01:09<00:42, 3.04s/batch]Cost: 2.321150779724121
Epoch 1/5:	63%	22/35 [01:13<00:43, 3.37s/batch]Cost: 2.2311508655548096
Epoch 1/5:	66%	23/35 [01:16<00:40, 3.38s/batch]Cost: 2.2411508560180664
Epoch 1/5:	69%	24/35 [01:19<00:35, 3.26s/batch]Cost: 2.2911508083343506
Epoch 1/5:	71%	25/35 [01:22<00:31, 3.16s/batch]Cost: 2.3111507892608643
Epoch 1/5:	74%	26/35 [01:27<00:31, 3.46s/batch]Cost: 2.271150827407837
Epoch 1/5:	77%	27/35 [01:31<00:29, 3.71s/batch]Cost: 2.331150770187378
Epoch 1/5:	80%	28/35 [01:34<00:24, 3.46s/batch]Cost: 2.2911508083343506
Epoch 1/5:	83%	29/35 [01:37<00:19, 3.28s/batch]Cost: 2.211150884628296
Epoch 1/5:	86%	30/35 [01:39<00:15, 3.17s/batch]Cost: 2.3011507987976074
Epoch 1/5:	89%	31/35 [01:43<00:13, 3.43s/batch]Cost: 2.201150894165039
Epoch 1/5:	91%	32/35 [01:47<00:10, 3.40s/batch]Cost: 2.321150779724121
Epoch 1/5:	94%	33/35 [01:50<00:06, 3.22s/batch]Cost: 2.3411507606506348
Epoch 1/5:	97%	34/35 [01:52<00:03, 3.10s/batch]Cost: 2.2211508750915527

Cost: 2.2611501216888428
Epoch: 0 Cost: 80.31078362464905

Epoch 5/5:	3%	1/35 [00:04<02:21, 4.17s/batch]Cost: 2.331150770187378
Epoch 5/5:	6%	2/35 [00:08<02:16, 4.13s/batch]Cost: 2.26115083694458
Epoch 5/5:	9%	3/35 [00:11<01:52, 3.51s/batch]Cost: 2.2811508178710938
Epoch 5/5:	11%	4/35 [00:13<01:39, 3.20s/batch]Cost: 2.3011507987976074
Epoch 5/5:	14%	5/35 [00:16<01:31, 3.05s/batch]Cost: 2.331150770187378
Epoch 5/5:	17%	6/35 [00:19<01:30, 3.13s/batch]Cost: 2.26115083694458
Epoch 5/5:	20%	7/35 [00:23<01:35, 3.42s/batch]Cost: 2.271150827407837
Epoch 5/5:	23%	8/35 [00:26<01:26, 3.21s/batch]Cost: 2.2911508083343506
Epoch 5/5:	26%	9/35 [00:29<01:19, 3.07s/batch]Cost: 2.2811508178710938
Epoch 5/5:	29%	10/35 [00:32<01:14, 2.97s/batch]Cost: 2.2811508178710938
Epoch 5/5:	31%	11/35 [00:35<01:14, 3.10s/batch]Cost: 2.321150779724121
Epoch 5/5:	34%	12/35 [00:39<01:15, 3.30s/batch]Cost: 2.3511507511138916
Epoch 5/5:	37%	13/35 [00:42<01:08, 3.13s/batch]Cost: 2.2511508464813232
Epoch 5/5:	40%	14/35 [00:44<01:03, 3.01s/batch]Cost: 2.3111507892608643
Epoch 5/5:	43%	15/35 [00:47<00:58, 2.93s/batch]Cost: 2.2811508178710938
Epoch 5/5:	46%	16/35 [00:51<00:59, 3.11s/batch]Cost: 2.3511507511138916
Epoch 5/5:	49%	17/35 [00:54<00:58, 3.26s/batch]Cost: 2.2311508655548096
Epoch 5/5:	51%	18/35 [00:57<00:52, 3.10s/batch]Cost: 2.3411507606506348
Epoch 5/5:	54%	19/35 [01:00<00:47, 2.99s/batch]Cost: 2.331150770187378
Epoch 5/5:	57%	20/35 [01:02<00:43, 2.93s/batch]Cost: 2.331150770187378
Epoch 5/5:	60%	21/35 [01:06<00:44, 3.17s/batch]Cost: 2.321150779724121
Epoch 5/5:	63%	22/35 [01:10<00:42, 3.27s/batch]Cost: 2.2311508655548096
Epoch 5/5:	66%	23/35 [01:12<00:37, 3.13s/batch]Cost: 2.2411508560180664
Epoch 5/5:	69%	24/35 [01:15<00:33, 3.01s/batch]Cost: 2.2911508083343506
Epoch 5/5:	71%	25/35 [01:18<00:29, 2.94s/batch]Cost: 2.3111507892608643
Epoch 5/5:	74%	26/35 [01:22<00:29, 3.22s/batch]Cost: 2.271150827407837
Epoch 5/5:	77%	27/35 [01:25<00:25, 3.25s/batch]Cost: 2.331150770187378
Epoch 5/5:	80%	28/35 [01:28<00:21, 3.09s/batch]Cost: 2.2911508083343506
Epoch 5/5:	83%	29/35 [01:31<00:17, 3.00s/batch]Cost: 2.211150884628296
Epoch 5/5:	86%	30/35 [01:33<00:14, 2.93s/batch]Cost: 2.3011507987976074
Epoch 5/5:	89%	31/35 [01:37<00:13, 3.28s/batch]Cost: 2.201150894165039
Epoch 5/5:	91%	32/35 [01:41<00:09, 3.25s/batch]Cost: 2.321150779724121
Epoch 5/5:	94%	33/35 [01:43<00:06, 3.11s/batch]Cost: 2.3411507606506348
Epoch 5/5:	97%	34/35 [01:46<00:03, 3.01s/batch]Cost: 2.2211508750915527

Epoch: 4 Cost: 80.14027762413025
Cost: 2.2611501216888428

```
# Initialize objects for counting correct/total
correct = 0

# pass each batch into the network
outputs = net(images)

# the class with the maximum score is what will be the predicted class
_, predicted = torch.max(outputs.data, 1)

# add size of the current batch
total += labels.size(0)

# add the number of correct predictions in the current batch
correct += (predicted == labels).sum().item()

# Calculate and print the proportion correct
print(correct/total)
```

Run cell (⌘/Ctrl+Enter)
cell executed since last change

executed by Sam Rudenberg
3:55 PM (42 minutes ago)
executed in 48.781s

0.17063142437591777

```
# Convert Pandas Series to NumPy array
test_y_array = test_y.to_numpy()

# Convert NumPy array to PyTorch tensor
test_y_tensor = torch.Tensor(test_y_array)

### Convert to numpy array then reshape
test_unflattened = test_X.reshape(len(test_y), 240, 140, 1)

## Convert test images into a tensor
test_tensor = torch.from_numpy(test_unflattened)

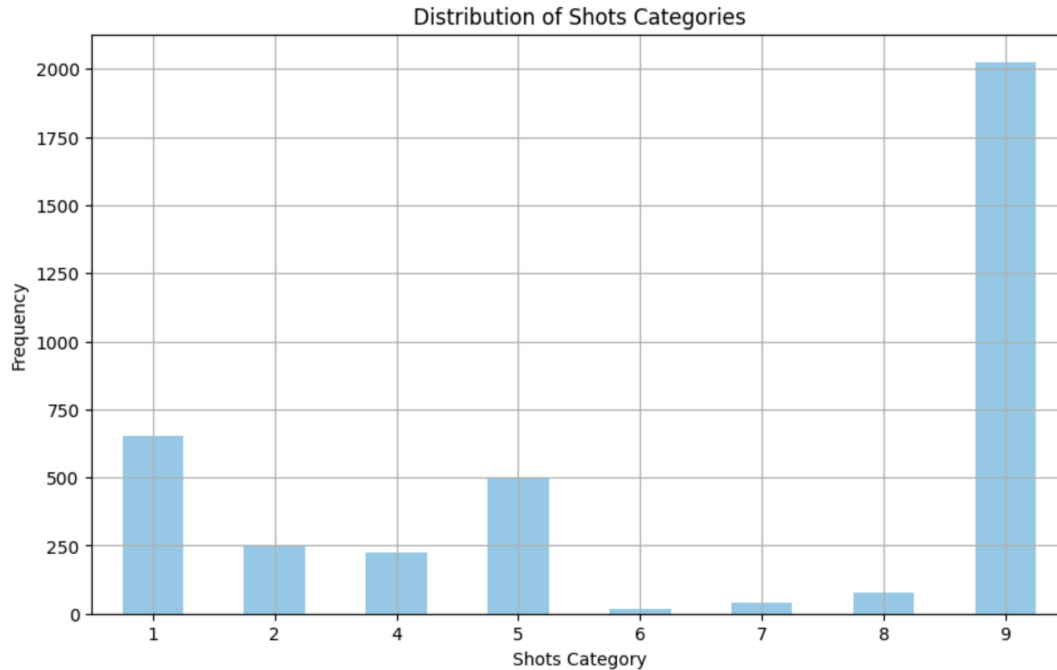
## Combine X and y tensors into a TensorDataset and DataLoader
test_loader = DataLoader(TensorDataset(test_tensor.type(torch.FloatTensor), test_y_tensor.type(torch.LongTensor)), batch_size=bsize)

## Repeat evaluation loop using the test data
correct = 0
total = 0
with torch.no_grad():
    for data in test_loader:
        images, labels = data
        outputs = net(images)
        _, predicted = torch.max(outputs.data, 1)
        total += labels.size(0)
        correct += (predicted == labels).sum().item()
print(correct/total)
```

0.18997361477572558

Above is the accuracy on the training and test sets.

Taking a deeper look at our data, there are class imbalances. The majority of frames in each clip are standby, and each clip always begins with a serve, however, the other shots appear considerably less, if at all. If we were to repeat this project, we would attempt to find videos to equalize our categories. Below is a distribution of the types of classifications across all clips.



Discussion

The models we embarked on doing were definitely very ambitious and required access to a more powerful pc that could be running the code in the background. Ideally we should have found a suitable pre-trained model to help us “deload” some of the pressure in the RAM. The big challenge involved taking into account the temporal aspect of the data since the images were labeled to be treated together. As expected, running a model that did not take time into account does not have an optimal accuracy, randomly guessing “standby” would provide us with 50% accuracy.

References

Miller, R. E. (n.d.). Machine Learning - Supervised Learning: Introduction to Regression [Webpage]. Retrieved from <https://remiller1450.github.io/MLs24.html>