Student Name: Mohammad Mohammad Beigi - Pantea Amoie
Student ID: 99102189 - 400101656
Subject: ML Privacy

**Introduction to Machine Learning - Dr. R. Amiri**
Final Project

# 1 Machine Unlearning

# 2 Private Training Attack

## 2.1 Base Model Learning

```
Train: X=(60000, 28, 28), y=(60000,)
Test: X=(10000, 28, 28), y=(10000,)
Model: "sequential_2"
```

| Layer (type) | Output Shape | Param # |
|---|---|---|
| flatten_2 (Flatten) | (None, 784) | 0 |
| dense_4 (Dense) | (None, 128) | 100,480 |
| dense_5 (Dense) | (None, 10) | 1,290 |

```
Total params: 101,770 (397.54 KB)
Trainable params: 101,770 (397.54 KB)
Non-trainable params: 0 (0.00 B)
Epoch 1/5
94/94 ──────────────── 2s 4ms/step - accuracy: 0.5073 - loss: 82.7176
Epoch 2/5
94/94 ──────────────── 1s 3ms/step - accuracy: 0.7224 - loss: 13.4081
Epoch 3/5
94/94 ──────────────── 0s 3ms/step - accuracy: 0.7444 - loss: 8.7707
Epoch 4/5
94/94 ──────────────── 1s 3ms/step - accuracy: 0.7929 - loss: 5.0982
Epoch 5/5
94/94 ──────────────── 1s 3ms/step - accuracy: 0.7845 - loss: 1.9695
63/63 ──────────────── 0s 2ms/step - accuracy: 0.6379 - loss: 3.7948

Model  0  test accuracy: 0.6439999938011169
```

Figure 1:

## 2.2 Base Model Shadow Learning and attack

```
Model  0  test accuracy: 0.6439999938011169
Epoch 1/5
94/94 ──────────────── 1s 3ms/step - accuracy: 0.4784 - loss: 90.4529
Epoch 2/5
94/94 ──────────────── 1s 3ms/step - accuracy: 0.7028 - loss: 16.0662
Epoch 3/5
94/94 ──────────────── 1s 3ms/step - accuracy: 0.7385 - loss: 9.1251
Epoch 4/5
94/94 ──────────────── 0s 3ms/step - accuracy: 0.7479 - loss: 7.1349
Epoch 5/5
94/94 ──────────────── 0s 3ms/step - accuracy: 0.7733 - loss: 5.2279
63/63 ──────────────── 0s 2ms/step - accuracy: 0.6955 - loss: 7.1460

Model  0  test accuracy: 0.6790000200271606
Epoch 1/5
94/94 ──────────────── 1s 3ms/step - accuracy: 0.4726 - loss: 82.6930
Epoch 2/5
94/94 ──────────────── 1s 3ms/step - accuracy: 0.6843 - loss: 16.4903
Epoch 3/5
94/94 ──────────────── 0s 3ms/step - accuracy: 0.7010 - loss: 5.9780
Epoch 4/5
94/94 ──────────────── 1s 5ms/step - accuracy: 0.7113 - loss: 1.6879
Epoch 5/5
94/94 ──────────────── 1s 5ms/step - accuracy: 0.7634 - loss: 1.2455
63/63 ──────────────── 0s 3ms/step - accuracy: 0.7023 - loss: 2.9090

Model  1  test accuracy: 0.7080000042915344
Training svm model :  0
[LibSVM]SVM model  0 score :  0.6102150537634409
```

Figure 2: Number of shadow models = 2

## 2.3 Private Model Learning

```
Train: X=(60000, 28, 28), y=(60000,)
Test: X=(10000, 28, 28), y=(10000,)
Model: "sequential_4"

_____
 Layer (type)                Output Shape              Param #
=================================================================
 flatten_4 (Flatten)         (None, 784)               0

 dense_8 (Dense)             (None, 128)               100480

 dropout_6 (Dropout)         (None, 128)               0

 dense_9 (Dense)             (None, 10)                1290

=================================================================
Total params: 101770 (397.54 KB)
Trainable params: 101770 (397.54 KB)
Non-trainable params: 0 (0.00 Byte)
_____

Epoch 1/5
94/94 [==============================] - 1s 4ms/step - loss: 24.4152 - accuracy: 0.3330
Epoch 2/5
94/94 [==============================] - 0s 3ms/step - loss: 2.3659 - accuracy: 0.4293
Epoch 3/5
94/94 [==============================] - 0s 3ms/step - loss: 1.9294 - accuracy: 0.4753
Epoch 4/5
94/94 [==============================] - 0s 4ms/step - loss: 1.7721 - accuracy: 0.5003
Epoch 5/5
94/94 [==============================] - 0s 3ms/step - loss: 1.7717 - accuracy: 0.4990
63/63 [==============================] - 0s 2ms/step - loss: 1.8605 - accuracy: 0.5755

Model  0  test accuracy: 0.5755000114440918
```

Figure 3: Private Trained Model by adding DropOut

## 2.4 Private Model Shadow Learning and attack

```
Model  0  test accuracy: 0.5755000114440918
Epoch 1/5
94/94 [==============================] - 1s 3ms/step - loss: 28.6290 - accuracy: 0.3900
Epoch 2/5
94/94 [==============================] - 0s 3ms/step - loss: 2.2916 - accuracy: 0.3817
Epoch 3/5
94/94 [==============================] - 0s 3ms/step - loss: 1.9308 - accuracy: 0.4120
Epoch 4/5
94/94 [==============================] - 0s 3ms/step - loss: 1.7087 - accuracy: 0.4530
Epoch 5/5
94/94 [==============================] - 0s 4ms/step - loss: 1.6937 - accuracy: 0.4740
63/63 [==============================] - 0s 2ms/step - loss: 1.3593 - accuracy: 0.6210

 Model  0  test accuracy: 0.6209999918937683
Epoch 1/5
94/94 [==============================] - 1s 3ms/step - loss: 24.4149 - accuracy: 0.3767
Epoch 2/5
94/94 [==============================] - 0s 3ms/step - loss: 2.3103 - accuracy: 0.3637
Epoch 3/5
94/94 [==============================] - 0s 4ms/step - loss: 1.9929 - accuracy: 0.4233
Epoch 4/5
94/94 [==============================] - 1s 6ms/step - loss: 1.8375 - accuracy: 0.4540
Epoch 5/5
94/94 [==============================] - 0s 5ms/step - loss: 1.7843 - accuracy: 0.4907
63/63 [==============================] - 0s 3ms/step - loss: 1.3897 - accuracy: 0.5945

 Model  1  test accuracy: 0.5945000052452087
Training svm model :  0
[LibSVM]SVM model  0  score :  0.48118279569892475
```

Figure 4: Number of shadow models = 2

## 2.5 Base Model Increasing Shadow Models

```
Model: "sequential_5"
```

| Layer (type) | Output Shape | Param # |
| --- | --- | --- |
| flatten_5 (Flatten) | (None, 784) | 0 |
| dense_10 (Dense) | (None, 128) | 100,480 |
| dropout_1 (Dropout) | (None, 128) | 0 |
| dense_11 (Dense) | (None, 10) | 1,290 |

```
Total params: 101,770 (397.54 KB)
Trainable params: 101,770 (397.54 KB)
Non-trainable params: 0 (0.00 B)
Epoch 1/5
94/94 ───────────────── 2s 8ms/step - accuracy: 0.4059 - loss: 85.0269
Epoch 2/5
94/94 ───────────────── 4s 35ms/step - accuracy: 0.4181 - loss: 2.2720
Epoch 3/5
94/94 ───────────────── 2s 3ms/step - accuracy: 0.4721 - loss: 1.8457
Epoch 4/5
94/94 ───────────────── 1s 3ms/step - accuracy: 0.4717 - loss: 2.0073
Epoch 5/5
94/94 ───────────────── 1s 14ms/step - accuracy: 0.5214 - loss: 1.6442
63/63 ───────────────── 1s 8ms/step - accuracy: 0.6274 - loss: 1.5189

 Model  0  test accuracy: 0.6119999885559082
Epoch 1/5
94/94 ───────────────── 2s 6ms/step - accuracy: 0.3803 - loss: 74.3437
Epoch 2/5
94/94 ───────────────── 2s 23ms/step - accuracy: 0.3819 - loss: 2.2994
Epoch 3/5
94/94 ───────────────── 0s 4ms/step - accuracy: 0.4140 - loss: 1.8617
Epoch 4/5
94/94 ───────────────── 1s 12ms/step - accuracy: 0.4580 - loss: 1.7917
Epoch 5/5
94/94 ───────────────── 1s 4ms/step - accuracy: 0.4706 - loss: 1.5689
63/63 ───────────────── 0s 2ms/step - accuracy: 0.6694 - loss: 1.3016

 Model  0  test accuracy: 0.6539999842643738
```

```
Model  0  test accuracy: 0.6539999842643738
Epoch 1/5
94/94 ───────────────── 2s 10ms/step - accuracy: 0.3444 - loss: 78.2884
Epoch 2/5
94/94 ───────────────── 1s 6ms/step - accuracy: 0.3500 - loss: 2.5628
Epoch 3/5
94/94 ───────────────── 0s 4ms/step - accuracy: 0.3818 - loss: 2.1204
Epoch 4/5
94/94 ───────────────── 0s 3ms/step - accuracy: 0.4212 - loss: 1.8821
Epoch 5/5
94/94 ───────────────── 1s 13ms/step - accuracy: 0.4528 - loss: 1.8724
63/63 ───────────────── 1s 6ms/step - accuracy: 0.5746 - loss: 1.5327

 Model  1  test accuracy: 0.5590000152587891
Epoch 1/5
94/94 ───────────────── 3s 16ms/step - accuracy: 0.3453 - loss: 68.9072
Epoch 2/5
94/94 ───────────────── 1s 6ms/step - accuracy: 0.4100 - loss: 2.5006
Epoch 3/5
94/94 ───────────────── 2s 5ms/step - accuracy: 0.4345 - loss: 2.1591
Epoch 4/5
94/94 ───────────────── 0s 3ms/step - accuracy: 0.4741 - loss: 1.7842
Epoch 5/5
94/94 ───────────────── 2s 21ms/step - accuracy: 0.4934 - loss: 1.6147
63/63 ───────────────── 0s 2ms/step - accuracy: 0.6377 - loss: 1.9835

 Model  2  test accuracy: 0.6305000185966492
Epoch 1/5
94/94 ───────────────── 3s 20ms/step - accuracy: 0.3645 - loss: 72.3915
Epoch 2/5
94/94 ───────────────── 1s 6ms/step - accuracy: 0.4231 - loss: 2.5470
Epoch 3/5
94/94 ───────────────── 2s 12ms/step - accuracy: 0.4441 - loss: 2.0397
Epoch 4/5
94/94 ───────────────── 0s 4ms/step - accuracy: 0.4795 - loss: 1.7557
Epoch 5/5
94/94 ───────────────── 1s 4ms/step - accuracy: 0.4982 - loss: 1.8641
63/63 ───────────────── 1s 15ms/step - accuracy: 0.6661 - loss: 1.3550

 Model  3  test accuracy: 0.6614999771118164
Training svm model :  0
[LibSVM]SVM model  0  score :  0.667632665465435
```

Figure 5: Number of shadow models = 4

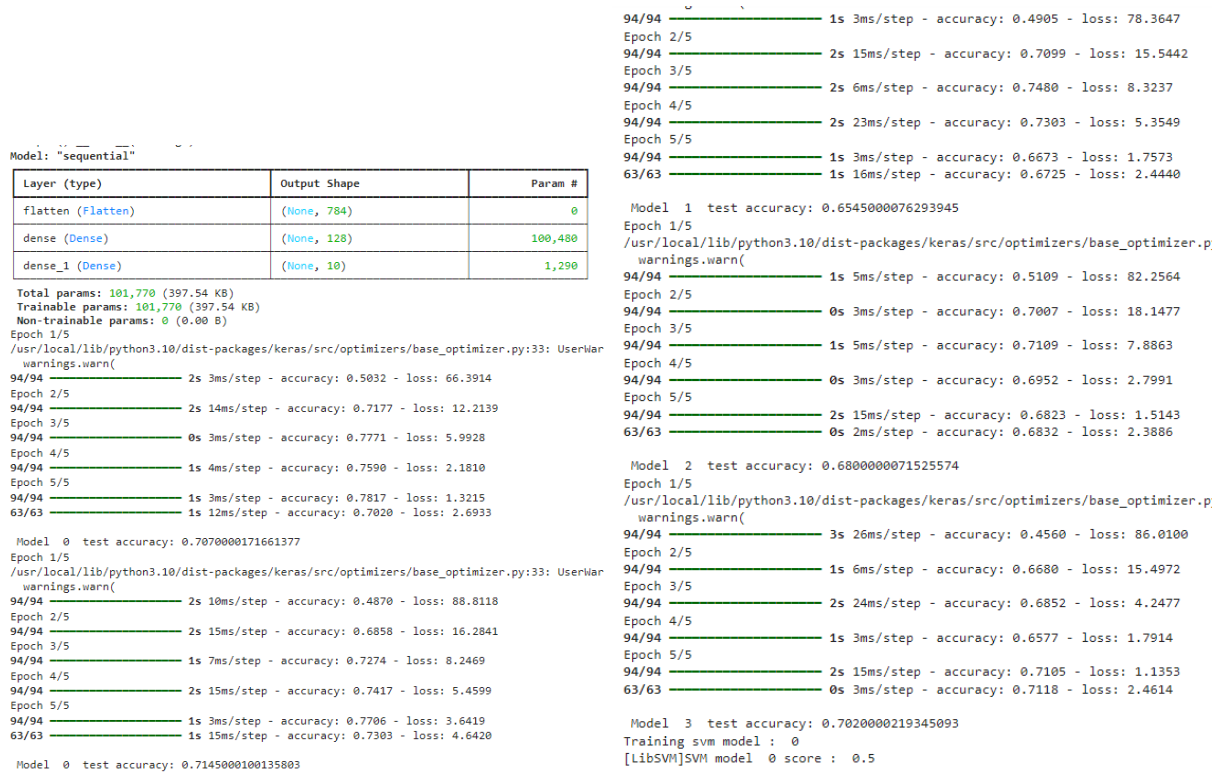## 2.6 Safe Model Increasing Shadow Models



Figure 6: Number of shadow models = 4

# 3 Membership Inference Attack

We trained an attacker with 2 shadow models similar to Question 4, but using the new model.

```python
from torchvision import models
import torch
import torchvision
import torchvision.transforms as transforms
from torch.utils.data import DataLoader
import torch.nn as nn
import torch.nn.functional as F
import torch.optim as optim
import numpy as np
import matplotlib.pyplot as plt


class CIFAR10Classifier(nn.Module):
    def __init__(self):
        super(CIFAR10Classifier, self).__init__()
        self.conv1 = nn.Conv2d(3, 16, 3, 1)
        self.conv2 = nn.Conv2d(16, 32, 3, 1)
        self.dropout1 = nn.Dropout2d(0.25)
        self.dropout2 = nn.Dropout2d(0.5)
        self.fc1 = nn.Linear(6272, 64)
        self.fc2 = nn.Linear(64, 10)

    def forward(self, x):
        x = self.conv1(x)
        x = F.relu(x)
```

```
26    x = self.conv2(x)
27    x = F.relu(x)
28    x = F.max_pool2d(x, 2)
29    x = self.dropout1(x)
30    x = torch.flatten(x, 1)
31    x = self.fc1(x)
32    x = F.relu(x)
33    x = self.dropout2(x)
34    x = self.fc2(x)
35    return x
36
37    import tensorflow as tf
38    import tensorflow.keras
39    from tensorflow.keras.utils import to_categorical
40    import numpy as np
41    from sklearn.utils import resample
42    from sklearn.metrics import ConfusionMatrixDisplay
43    import matplotlib.pyplot as plt
44    from matplotlib import pyplot
45
46
47
48    def load_fashion_mnist_dataset():
49    from keras.datasets import fashion_mnist
50
51    # load dataset
52    (trainX, trainY), (testX, testY) = fashion_mnist.load_data()
53    # summarize loaded dataset
54    print('Train: X=%s, y=%s' % (trainX.shape, trainY.shape))
55    print('Test: X=%s, y=%s' % (testX.shape, testY.shape))
56
57    # reshape dataset to have a single channel
58    trainX = trainX.reshape((trainX.shape[0], 28, 28, 1))
59    testX = testX.reshape((testX.shape[0], 28, 28, 1))
60
61    # one hot encode target values
62    trainY = to_categorical(trainY)
63    testY = to_categorical(testY)
64
65    return (trainX, trainY), (testX, testY)
66
67    # Set hyper-parameters for our training of neural networks
68    LEARNING_RATE = 0.001
69    EPOCH = 5
70    VERBOSE = 1
71
72    TRAINING_SIZE = 3000
73    TEST_SIZE = 2000
74
75    # No of target models
76    NUM_TARGET = 1
77    # No of shadow models
78    NUM_SHADOW = 2
79
80    # Label value "in" for records present in training data of shadow models
81    IN = 1
82    # Label value "out" for records not present in training data of shadow models
83    OUT = 0
84
85    def sample_data(train_data, test_data, num_sets):
86    (x_train, y_train), (x_test, y_test) = train_data, test_data
```

```python
87     new_x_train, new_y_train = [], []
88     new_x_test, new_y_test = [], []
89     for i in range(num_sets):
90         x_temp, y_temp = resample(x_train, y_train, n_samples=TRAINING_SIZE, random_state=0)
91         new_x_train.append(x_temp)
92         new_y_train.append(y_temp)
93         x_temp, y_temp = resample(x_test, y_test, n_samples=TEST_SIZE, random_state=0)
94         new_x_test.append(x_temp)
95         new_y_test.append(y_temp)
96     return (new_x_train, new_y_train), (new_x_test, new_y_test)
97
98     def get_attack_dataset(models, train_data, test_data, num_models, data_size):
99         # generate dataset for the attack model
100        (x_train, y_train), (x_test, y_test) = train_data, test_data
101        # set number of classes for the attack model
102        num_classes = 10
103        x_data, y_data = [[] for _ in range(num_classes)], [[] for _ in range(num_classes)]
104        for i in range(num_models):
105            # IN data
106            x_temp, y_temp = resample(x_train[i], y_train[i], n_samples=data_size, random_state=0)
107            for j in range(data_size):
108                y_idx = np.argmax(y_temp[j])
109                x_data[y_idx].append(models[i].predict(x_temp[j:j+1], verbose=0)[0])
110                y_data[y_idx].append(IN)
111
112            # OUT data
113            x_temp, y_temp = resample(x_test[i], y_test[i], n_samples=data_size, random_state=0)
114            for j in range(data_size):
115                y_idx = np.argmax(y_temp[j])
116                x_data[y_idx].append(models[i].predict(x_temp[j:j+1], verbose=0)[0])
117                y_data[y_idx].append(OUT)
118
119        return x_data, y_data
120
121    def build_fcnn_model_fashion_mnist():
122        model = tf.keras.models.Sequential([
123            tf.keras.layers.Flatten(input_shape=(28, 28)),
124            tf.keras.layers.Dense(128, activation='relu'),
125            tf.keras.layers.Dense(10, activation='softmax')
126        ])
127        model.summary()
128        return model
129
130    def get_trained_keras_models(keras_model, train_data, test_data, num_models):
131        (x_train, y_train), (x_test, y_test) = train_data, test_data
132        models = []
133        for i in range(num_models):
134            models.append(tf.keras.models.clone_model(keras_model))
135            rms = tf.keras.optimizers.RMSprop(learning_rate=LEARNING_RATE, decay=1e-7)
136            models[i].compile(loss='categorical_crossentropy', optimizer=rms, metrics=['accuracy'])
137            models[i].fit(x_train[i], y_train[i], batch_size=32, epochs=EPOCH, verbose=VERBOSE,
                    shuffle=True)
138            score = models[i].evaluate(x_test[i], y_test[i], verbose=VERBOSE)
139            print('\n', 'Model ', i, ' test accuracy:', score[1])
140        return models
141
142    def get_trained_svm_models(train_data, test_data, num_models=1):
143        from sklearn import svm
144        (x_train, y_train), (x_test, y_test) = train_data, test_data
145        models = []
146        for i in range(num_models):
```

6

```python
147    print('Training svm model : ', i)
148    models.append(svm.SVC(gamma='scale', kernel='linear', verbose=VERBOSE))
149    models[i].fit(x_train[i], y_train[i])
150    score = models[i].score(x_test[i], y_test[i])
151    print('SVM model ', i, 'score : ', score)
152    return models
153
154    def membership_attack():
155    # load the pre-shuffled train and test data
156    (x_train, y_train), (x_test, y_test) = load_fashion_mnist_dataset()
157
158    # split the data for each model
159    target_train = (x_train[:TRAINING_SIZE*NUM_TARGET], y_train[:TRAINING_SIZE*NUM_TARGET])
160    target_test = (x_test[:TEST_SIZE*NUM_TARGET], y_test[:TEST_SIZE*NUM_TARGET])
161    target_train_data, target_test_data = sample_data(target_train, target_test, NUM_TARGET)
162
163    shadow_train = (x_train[TRAINING_SIZE*NUM_TARGET:], y_train[TRAINING_SIZE*NUM_TARGET:])
164    shadow_test = (x_test[TEST_SIZE*NUM_TARGET:], y_test[TEST_SIZE*NUM_TARGET:])
165    shadow_train_data, shadow_test_data = sample_data(shadow_train, shadow_test, NUM_SHADOW)
166
167    cnn_model = build_fcnn_model_fashion_mnist()
168
169    # compile the target model
170    target_models = get_trained_keras_models(cnn_model, target_train_data, target_test_data,
         NUM_TARGET)
171    # compile the shadow models
172    shadow_models = get_trained_keras_models(cnn_model, shadow_train_data, shadow_test_data,
         NUM_SHADOW)
173
174    # get train data for the attack model
175    attack_train = get_attack_dataset(shadow_models, shadow_train_data, shadow_test_data,
         NUM_SHADOW, TEST_SIZE)
176    # get test data for the attack model
177    attack_test = get_attack_dataset(target_models, target_train_data, target_test_data,
         NUM_TARGET, TEST_SIZE)
178
179    # training the attack model
180    #attack_model = get_trained_svm_models(attack_train, attack_test)
181    return attack_train, attack_test
182    NUM_SHADOW = 2
183    membership_attack()
184
185
186    import torch
187    import torch.nn as nn
188    import torch.optim as optim
189    from torchvision.datasets import CIFAR10
190    from torchvision import transforms
191    from torch.utils.data import Subset, DataLoader, TensorDataset
192    from sklearn.metrics import confusion_matrix, precision_score, recall_score ,f1_score
193    from sklearn.linear_model import LogisticRegression
194
195    device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
196
197    model = CIFAR10Classifier()
198    state_dict = torch.load("model_state_dict.pth", map_location=device)
199    new_state_dict = {key.replace('_module.', ''): value for key, value in state_dict.items()}
200    model.load_state_dict(new_state_dict)
201    model.to(device)
202    model.eval()
203
```

```
204    transform = transforms.Compose([
205    transforms.ToTensor(),
206    transforms.Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5))
207    ])
208
209    DATA_ROOT = '../cifar10'
210    BATCH_SIZE = 64
211
212    # Load the indices from list.txt
213    indices_file = 'list.txt' ############
214    with open(indices_file, 'r') as f:
215    indices = [int(line.strip()) for line in f]
216
217    full_train_dataset = CIFAR10(root=DATA_ROOT, train=True, download=True, transform=
         transform)
218    test_dataset = CIFAR10(root=DATA_ROOT, train=False, download=True, transform=transform)
219
220    train_indices_set = set(indices)
221    all_indices = set(range(len(full_train_dataset)))
222    other_indices = list(all_indices - train_indices_set)
223
224    train_dataset = Subset(full_train_dataset, indices[:len(indices)//2])  ###########
225    other_dataset = Subset(full_train_dataset, other_indices)
226
227    # Create data loaders
228    train_loader = DataLoader(train_dataset, batch_size=BATCH_SIZE, shuffle=False)
229    other_loader = DataLoader(other_dataset, batch_size=BATCH_SIZE, shuffle=False)
230    test_loader = DataLoader(test_dataset, batch_size=BATCH_SIZE, shuffle=False)
231
232    # Create labels
233    train_labels = torch.ones(len(train_dataset)).to(device)
234    other_labels = torch.zeros(len(other_dataset)).to(device)
235    test_labels = torch.zeros(len(test_dataset)).to(device)
236    ###################################
237    #if you have an attacker model for each class, modify the above code.
238    ###################################
239
240    def extract_features(model, dataloader):
241    model.eval()
242    features = []
243    with torch.no_grad():
244    for data in dataloader:
245    inputs, _ = data
246    inputs = inputs.to(device)
247    outputs = model(inputs)
248    features.append(outputs)
249    return torch.cat(features).to(device)
250
251    train_features = extract_features(model, train_loader)
252    other_features = extract_features(model, other_loader)
253    test_features = extract_features(model, test_loader)
254
255
256    combined_features = torch.cat((train_features, other_features, test_features))
257    combined_labels = torch.cat((train_labels, other_labels, test_labels))
258
259
260    new_dataset = TensorDataset(combined_features, combined_labels)
261    new_loader = DataLoader(new_dataset, batch_size=BATCH_SIZE, shuffle=True)
262
263    #load your attacker model
```

```
264
265    ##############################################
266
267    # Calculate training accuracy, confusion matrix, precision, and recall
268    binary_classifier.eval()
269    all_labels = []
270    all_predicted = []
271    correct = 0
272    total = 0
273    with torch.no_grad():
274    for features, labels in new_loader:
275    features, labels = features.to(device), labels.to(device)
276    outputs = get_trained_svm_models(features, attack_test).squeeze()
277    predicted = (outputs > 0.5).float()
278    total += labels.size(0)
279    correct += (predicted == labels).sum().item()
280    all_labels.extend(labels.cpu().numpy())
281    all_predicted.extend(predicted.cpu().numpy())
282
283    accuracy = correct / total
284
285
286    cm = confusion_matrix(all_labels, all_predicted)
287    precision = precision_score(all_labels, all_predicted)
288    recall = recall_score(all_labels, all_predicted)
289    f1 = f1_score(all_labels, all_predicted)
290
291    print(f'Confusion Matrix:\n{cm}')
292    print(f"Precision: {precision:.4f}")
293    print(f"Recall: {recall:.4f}")
294    print(f"F1 Score: {f1:.4f}")
295    print(f'Training Accuracy: {accuracy:.4f}')
```

```
Confusion Matrix:
[[ 9452  8766]
 [ 8361 13421]]
Precision: 0.6049
Recall: 0.6162
F1 Score: 0.6105
Training Accuracy: 0.5718
```

Figure 7:

# PyTorch and CIFAR-10 Classifier

**Explanation:**

1. **Imports**:

   - `torchvision`, `torch`, `transforms`: For handling data transformations and model-related functionalities.
   - `nn`, `F`: For defining neural network layers and activation functions.
   - `optim`: For optimization algorithms.
   - `numpy`, `matplotlib.pyplot`: For numerical operations and plotting.

2. **CIFAR10Classifier Class**:

- Inherits from `nn.Module`.
- **Initialization (\_\_init\_\_)**: Defines the structure of the neural network.
  - Convolutional layers (`conv1`, `conv2`), dropout layers (`dropout1`, `dropout2`), and fully connected layers (`fc1`, `fc2`).
- **Forward Method**: Defines the forward pass through the network.
  - Applies convolution, activation (ReLU), pooling, dropout, and flattening operations sequentially.

# TensorFlow and CIFAR 10 Dataset

**Explanation:**

1. **Imports**:

- `tensorflow`, `keras.utils`: For building and managing neural networks.
- `numpy`, `sklearn.utils`, `sklearn.metrics`, `matplotlib.pyplot`: For data handling, resampling, and plotting.

2. **load\_fashion\_mnist\_dataset Function**:

- Loads the Fashion-MNIST dataset.
- Reshapes the dataset to have a single channel.
- Converts the labels to categorical format (one-hot encoding).
- Returns the training and testing datasets.

# Sampling and Training

**Explanation:**

1. **sample\_data Function**:

- Resamples training and testing datasets to create multiple sets.
- Helps in creating diverse training and testing sets for training models.

2. **get\_attack\_dataset Function**:

- Generates datasets for training attack models.
- Uses predictions from multiple models on resampled data to create in-distribution (IN) and out-of-distribution (OUT) samples.

# Neural Network and SVM Models

**Explanation:**

1. **build\_fcnn\_model\_fashion\_mnist Function**:

- Builds a fully connected neural network for Fashion-MNIST.
- Uses two dense layers with ReLU and softmax activations.

2. **get_trained_keras_models Function**:

   - Trains multiple instances of a given Keras model.
   - Uses RMSprop optimizer and evaluates each model.

3. **get_trained_svm_models Function**:

   - Trains multiple SVM models using the provided training data.
   - Evaluates each model and prints the accuracy.

# Membership Attack

**Explanation:**

1. **membership_attack Function**:

   - Implements a membership inference attack.
   - Loads and preprocesses the Fashion-MNIST dataset.
   - Splits data for training target and shadow models.
   - Builds and trains neural networks (both target and shadow models).
   - Creates datasets for training an attack model using predictions from shadow models.
   - Returns the attack training and testing datasets.

# Evaluating CIFAR-10 Classifier with Torch

**Explanation:**

1. **Device Setup**:

   - Checks for GPU availability and sets the device accordingly.

2. **Model Loading**:

   - Loads the pretrained CIFAR10Classifier model state.

3. **Data Preparation**:

   - Defines transformations for the CIFAR-10 dataset.
   - Loads CIFAR-10 dataset and splits it based on indices from a file.

4. **Feature Extraction**:

   - Extracts features from the trained model for training, other, and test datasets.

5. **Creating New Dataset**:

   - Combines extracted features and labels into a new dataset.
   - Creates a DataLoader for the new dataset.

6. **Evaluation**:

   - Evaluates the binary classifier on the new dataset.
   - Computes accuracy, confusion matrix, precision, recall, and F1 score.