# DATA PREPARATION

- Considering the CIFAR - 10 dataset
- Converting this dataset into binary dataset [which have label as 0 or 1]
- Converting this dataset into unbalalanced dataset by converting labels '0' as '0' and rest eveything label as '1'

```python
1 from torchvision.datasets import CIFAR10
2 import torchvision.transforms as transforms
3 import os
4 import pandas as pd
5 from torchvision.transforms import ToPILImage
6
7 transform = transforms.ToTensor()
8
9 train_set = CIFAR10(root='./data', train=True, download=Tru
10 test_set = CIFAR10(root='./data', train=False, download=Tru
11
```

```
100%|████████████| 170M/170M [00:15<00:00, 11.2MB/s]
```

```python
1 def label_adjustment(dataset):
2     results = []
3     for img, label in dataset:
4         if label == 0:
5             results.append([img, label])
6         if label == 1 or label == 2:
7             results.append([img, 1])
8     return results
```

```python
1 train_data = label_adjustment(train_set)
2 test_data  = label_adjustment(test_set)
```

```python
1 save_root = '/content/drive/MyDrive/cifar10_binary'
2 os.makedirs(save_root, exist_ok = True)
3 to_pil = ToPILImage()
4
5 def save_images_and_make_csv(data, split_name):
6     dir = os.path.join(save_root, split_name)
7     os.makedirs(dir, exist_ok = True)
8     rows = []
```

```
 9        for idx, (img_tensor, label) in enumerate(data):
10            img_path = os.path.join(dir, f'{idx}.png')
11            to_pil(img_tensor).save(img_path)
12            rows.append([img_path, label])
13
14        df = pd.DataFrame(rows, columns=["MD5HASH", "LABEL"])
15        df.to_csv(os.path.join(save_root, f"{split_name}.csv"),
16        print(f"{split_name}.csv saved with {len(rows)} entries
17
18 # Save both splits
19 save_images_and_make_csv(train_data, "train")
20 save_images_and_make_csv(test_data, "test")
21
```

```
train.csv saved with 15000 entries.
test.csv saved with 3000 entries.
```

```
 1 !ls /content/drive/MyDrive/cifar10_binary
```

```
test  test.csv  train  train.csv
```

## IMPORTING LIBRARIES AND SETTING CONGIURATION PARAMETERS

```
 1 from torch.utils.data import Dataset, DataLoader
 2 import torch
 3 import torch.nn as nn
 4 import torch.nn.functional as F
 5 from sklearn.metrics import classification_report
 6 from transformers import AutoImageProcessor, AutoModelForIm
 7 from torch.optim.lr_scheduler import StepLR
 8 # note: PIL stands for pillow; to install type "pip3 instal
 9 from PIL import Image
10 from datetime import datetime
11 from torchvision.transforms import Compose, Resize, RandomR
12 import pandas as pd
13 import numpy as np
14 import os, gc
```

```
2025-06-19 07:17:03.011678: E external/local_xla/xla/stream_executor/cuda/cud
WARNING: All log messages before absl::InitializeLog() is called are written
E0000 00:00:1750317423.209937      70 cuda_dnn.cc:8310] Unable to register cu
E0000 00:00:1750317423.267046      70 cuda_blas.cc:1418] Unable to register c
```

```python
1  # Check for CUDA and MPS availability, set the device accor
2  if torch.backends.mps.is_available():
3      device = torch.device("mps")
4      # setting environment variables, need to run training i
5      os.environ['PYTORCH_MPS_HIGH_WATERMARK_RATIO'] = '0.0'
6      os.environ['PYDEVD_DISABLE_FILE_VALIDATION'] = '1'
7      print("Using MPS as the device.")
8  else:
9      if torch.cuda.is_available():
10     # the syntax 'cuda:3' used to point a specific GPU from
11     # 'cuda' points to first GPU from the cluster of GPUs
12         device = torch.device("cuda")
13         print("Using CUDA as the device.")
14     else:
15         device = torch.device("cpu")
16         print("Using CPU as the device.")
```

Using CUDA as the device.

```python
1  TRAINING_DATA                              = '/content/drive/MyD
2  TESTING_DATA                               = '/content/drive/MyD
3  BATCH_SIZE                                 = 32#256
4  WORKERS                                    = 4
5  PIN_MEMORY                                 = True
6  MIXING                                     = True
7  MODEL_NAME                                 = "facebook/dinov2-ba
8  RESULTS                                    = 'results'
9  EPOCHS                                     = 4
10 BEST_MODEL                                 = None
11 PRETRAINING                                = False
12 LEARNING_RATE                              = 5e-5
13 L2_PENALTY                                 = 1e-5
14 GAMMA                                      = 0.1
15 STEPSIZE                                   = 3
16 SAVE_CHECKPOINTS                           = True
17 MIN_LOSS                                   = float('inf')
18 MODEL_SAVED                                = f'{RESULTS}/bestmod
19 THRESHOLD                                  = 0.5
20 OUTPUT_DIM                                 = 1
21
```

## PERFORMANCE METRIC [PRECISION & RECALL]

```
1 def calculate_classification_accuracy(loader, model):
2    model.eval()  # Set the model in evaluation mode
3    LABELS = []
4    PREDICTIONS = []
5
6    with torch.no_grad():
7        for images, labels in loader:
8            # Move to device and cast to float32
9            images, labels = images.to(device), labels.to(d
10           probabilities = model(images).squeeze()
11           # Predictions based on the threshold
12           prediction = torch.where(probabilities > THRESH
13           LABELS.extend(labels.tolist())
14           PREDICTIONS.extend(prediction.tolist())
15   return classification_report(LABELS, PREDICTIONS)
16
17
```

## DINO MODEL [WITH ADDITIONAL CLASSIFICATION HEAD]

```
1  class Network(nn.Module):
2     def __init__(self):
3         super(Network, self).__init__()
4
5         # taking processor for necessary substitions, if ne
6         self.processor = AutoImageProcessor.from_pretrained
7         self.model = Dinov2Model.from_pretrained(MODEL_NAME
8         self.pretrained_model_last_dim = self.model.layerno
9
10        # Additional classification head used for downstrea
11        self.cls_head = nn.Sequential(
12            nn.Linear(self.pretrained_model_last_dim, OUTPU
13            nn.Sigmoid(),
14        )
15
16        self.initialize_weights() # weights initialization
17
18    def initialize_weights(self):
19        torch.manual_seed(444)
20        for layer in self.cls_head:
21            if isinstance(layer, nn.Linear):
22                nn.init.kaiming_uniform_(layer.weight, non
23                nn.init.zeros_(layer.bias)
```

```
23                          ................................_(................)
24                          print(f"kaiming_uniform_ Initialization: {
25
26
27      def forward(self, x):
28          x = self.model(x).last_hidden_state[:, 0]
29          x = self.cls_head(x)
30          return x
31
32
```

```
1 model = Network()
2
3 model.to(device)
4
5
```

```
preprocessor_config.json:   0%|          | 0.00/436 [00:00<?, ?B/s]
Using a slow image processor as `use_fast` is unset and a slow processor was
config.json:   0%|          | 0.00/548 [00:00<?, ?B/s]
model.safetensors:   0%|          | 0.00/346M [00:00<?, ?B/s]
kaiming_uniform_ Initialization: Linear
Network(
  (model): Dinov2Model(
    (embeddings): Dinov2Embeddings(
      (patch_embeddings): Dinov2PatchEmbeddings(
        (projection): Conv2d(3, 768, kernel_size=(14, 14), stride=(14, 14))
      )
      (dropout): Dropout(p=0.0, inplace=False)
    )
    (encoder): Dinov2Encoder(
      (layer): ModuleList(
```

## DATA LOADER WITH UPSAMPLING

```
        (norm1): LayerNorm((768,), eps=1e-06, elementwise_affine=True)
        (attention): Dinov2Attention(
```

```python
 1
 2 class MD5HASHDataset(Dataset):
 3     def __init__(self, dataframe):
 4         self.dataframe = dataframe
 5         self.images = self.dataframe['MD5HASH'].values
 6         self.labels = self.dataframe['LABEL'].values
 7         self.processor = model.processor
 8         self.mean = self.processor.image_mean
 9         self.std = self.processor.image_std
10         self.interpolation = self.processor.resample
11
12         self.train_transform = Compose([
13             Resize(size = (32, 32)),
14             #Resize(size = (85, 550)),
15             #RandomResizedCrop(size = (224, 224),
16             #                  scale = (0.08, 1.0),
17             #                  ratio = (0.75, 1.3333),
18             #                  interpolation = self.interpo
19             #RandomHorizontalFlip(p = 0.5),
20             #ColorJitter(brightness = (0.6, 1.4),
21             #            contrast = (0.6, 1.4),
22             #            saturation = (0.6, 1.4)),
23             ToTensor(),
24             Normalize(mean = self.mean, std = self.std),
25         ])
26
27
28     def __len__(self):
29         return len(self.images)
30
31     def __getitem__(self, idx):
```

```
32          # Load the image from the file path
33          image_path = self.images[idx]
34          image = self.train_transform(Image.open(image_path)
35          # Get the label
36          label = torch.tensor(self.labels[idx], dtype=torch.
37
38          return image, label
39
```

```
1  def create_training_loader(data_csv = TRAINING_DATA, upsamp
2      # Load data
3      training_data = pd.read_csv(data_csv)
4      print(':::: TRAINING DATA DETAILS ::::')
5      print('- Number of Samples:', training_data.shape[0])
6      print('- LABEL DISTRIBUTION: \n',training_data['LABEL']
7
8      # Create dataset and dataloader
9      md5hash_dataset = MD5HASHDataset(training_data)
10     if upsampling:
11         # References:
12         # https://pytorch.org/docs/stable/data.html
13         # https://towardsdatascience.com/demystifying-pytor
14         from torch.utils.data import WeightedRandomSampler
15         classes_count = dict(training_data['LABEL'].value_c
16         sample_weights = [ 1 / classes_count[i] for i in tr
17         sampler = WeightedRandomSampler(weights = sample_we
18                                         num_samples = len(t
19                                         replacement = True)
20         data_loader = DataLoader(md5hash_dataset,
21                                  batch_size = BATCH_SIZE,
22                                  num_workers = WORKERS,
23                                  pin_memory = PIN_MEMORY,
24                                  shuffle = False,
25                                  sampler = sampler)
26     else:
27         data_loader = DataLoader(md5hash_dataset,
28                                  batch_size = BATCH_SIZE,
29                                  num_workers = WORKERS,
30                                  pin_memory = PIN_MEMORY,
31                                  shuffle = MIXING)
32
33     # Clean memory, :)
34     del training_data
35
```

```
36        return data_loader
37
38
```

```
1 data_loader = create_training_loader(upsampling = True)
2
```

```
::: TRAINING DATA DETAILS :::
- Number of Samples: 15000
- LABEL DISTRIBUTION:
 LABEL
1    10000
0     5000
Name: count, dtype: int64
```

## OPTIMIZER, SCHEDULER AND LOSS FUNCTION

```
1 # ## MODEL TRAINING
2
3 # Define the optimizer
4 # Idea borrowed from Research paper titled as "Improving Ge
5 if PRETRAINING:
6     optimizer = torch.optim.SGD(model.parameters(), lr = LE
7 else:
8     optimizer = torch.optim.Adam(model.parameters(), lr = L
9
10 # Define a learning rate scheduler
11 scheduler = StepLR(optimizer, step_size = STEPSIZE, gamma =
12 # Define the loss function: BCE
13 criterion = nn.BCELoss()
14
15
```

```
1 # directory creation
2 os.makedirs(RESULTS, exist_ok = True)
3 if SAVE_CHECKPOINTS:
4     CHECKPOINTDIR = f'{RESULTS}/checkpoints'
5     os.makedirs(CHECKPOINTDIR, exist_ok = True)
6
```

## MODEL FINE-TUNING TRAINING

```
1  # TRAINING LOOP
2  for epoch in range(EPOCHS):
3      print('-'*70)
4      # Define the total number of batches in the loader
5      total_loss = 0.0
6
7      # setting model stage to training
8      model.train()
9
10     for batch_idx, (images, labels) in enumerate(data_loade
11         # shifting to MPS
12         # Shift to MPS and then cast to float32
13         images, labels = images.to(device), labels.to(devic
14
15         # Forward pass
16         optimizer.zero_grad()  # Moved this line here to av
17
18         with torch.set_grad_enabled(True):
19             # Forward pass
20             outputs = model(images).squeeze()  # Squeeze to
21             loss = criterion(outputs, labels)
22
23             loss.backward()
24             optimizer.step()
25
26         total_loss += loss.item()
27         # Explicitly free up GPU memory
28         if torch.backends.mps.is_available():
29             torch.backends.mps.is_macos13_or_newer.cache_cl
30         if torch.cuda.is_available():
31             torch.cuda.empty_cache()
32         # Run garbage collector to free up CPU memory
33         gc.collect()
34
35     print(f"Epoch {epoch + 1}/{EPOCHS}, Loss: {total_loss /
36     print('TRAINING DATA')
37     print(f'- Performance: \n{calculate_classification_accu
38     # Update the learning rate
39     scheduler.step()
40
41
42     if SAVE_CHECKPOINTS:
43         timestamp = datetime.now().strftime("%Y%m%d%H%M%S")
44         checkpointmodel = '{}/epoch_{}_{}.pth'.format(CHECK
45         print('Saving checkpoint: ', checkpointmodel)
```

```
46            torch.save(model.state_dict(), checkpointmodel)
47
48      # Check if this epoch had the minimum loss
49      if total_loss < MIN_LOSS:
50          MIN_LOSS = total_loss
51          best_model = model.state_dict()
52          # Save the best model
53          if best_model is not None:
54              print('Saving Best Model: ', MODEL_SAVED)
55              torch.save(best_model, MODEL_SAVED)
56
57 ###############################
58
```

```
---------------------------------------------------------------
Epoch 1/4, Loss: 0.4652398204101301
TRAINING DATA
- Performance:
              precision    recall  f1-score   support

         0.0       0.88      0.90      0.89      7477
         1.0       0.90      0.87      0.89      7523

    accuracy                           0.89     15000
   macro avg       0.89      0.89      0.89     15000
weighted avg       0.89      0.89      0.89     15000

Saving checkpoint:  results/checkpoints/epoch_1_20250619072045.pth
Saving Best Model:  results/bestmodel.pth
---------------------------------------------------------------
Epoch 2/4, Loss: 0.18003429818366254
TRAINING DATA
- Performance:
              precision    recall  f1-score   support

         0.0       0.92      0.98      0.95      7466
         1.0       0.98      0.91      0.94      7534

    accuracy                           0.94     15000
   macro avg       0.95      0.94      0.94     15000
weighted avg       0.95      0.94      0.94     15000

Saving checkpoint:  results/checkpoints/epoch_2_20250619072407.pth
Saving Best Model:  results/bestmodel.pth
---------------------------------------------------------------
Epoch 3/4, Loss: 0.12192414922000312
TRAINING DATA
- Performance:
              precision    recall  f1-score   support

         0.0       0.97      0.95      0.96      7550
         1.0       0.95      0.97      0.96      7450

    accuracy                           0.96     15000
   macro avg       0.96      0.96      0.96     15000
```

```
weighted avg       0.96      0.96      0.96      15000

Saving checkpoint:  results/checkpoints/epoch_3_20250619072732.pth
Saving Best Model:  results/bestmodel.pth
------------------------------------------------------------------
Epoch 4/4, Loss: 0.05339157602428071
TRAINING DATA
- Performance:
              precision    recall  f1-score   support

         0.0       0.99      1.00      0.99      7483
         1.0       1.00      0.99      0.99      7517

    accuracy                           0.99     15000
   macro avg       0.99      0.99      0.99     15000
weighted avg       0.99      0.99      0.99     15000
```

```
1 model.load_state_dict(torch.load(MODEL_SAVED, weights_only
```

<All keys matched successfully>

## RESULTS ON ORIGINAL TRAINING AND TEST DATA

```
1 data_loader = create_training_loader() # without upsampling
2
3 # getting Best model performance on training, validation an
4 print('TRAINING DATA')
5 print(f'- Performance: \n{calculate_classification_accuracy
```

```
::: TRAINING DATA DETAILS :::
- Number of Samples: 15000
- LABEL DISTRIBUTION:
 LABEL
1    10000
0     5000
Name: count, dtype: int64
TRAINING DATA
- Performance:
              precision    recall  f1-score   support

         0.0       0.98      1.00      0.99      5000
         1.0       1.00      0.99      0.99     10000

    accuracy                           0.99     15000
   macro avg       0.99      0.99      0.99     15000
weighted avg       0.99      0.99      0.99     15000
```

```
1 testset_loader = create_training_loader(data_csv = TESTING_
2
3 print('TESTING DATA')
```

```
4 print(f'- Performance: \n{calculate_classification_accuracy
5
```

```
::: TRAINING DATA DETAILS :::
- Number of Samples: 3000
- LABEL DISTRIBUTION:
 LABEL
1    2000
0    1000
Name: count, dtype: int64
TESTING DATA
- Performance:
              precision    recall  f1-score   support

         0.0       0.90      0.88      0.89      1000
         1.0       0.94      0.95      0.95      2000

    accuracy                           0.93      3000
   macro avg       0.92      0.92      0.92      3000
weighted avg       0.93      0.93      0.93      3000
```