

The screenshot displays a Google Colab notebook titled 'Untitled0.ipynb'. The interface includes a top navigation bar with tabs for 'Untitled6.ipynb - Colab', 'destilbert - Colab', '(25) WhatsApp', 'Welcome To Colab - Colab', and 'Untitled0.ipynb - Colab'. Below this is a menu bar with 'File', 'Edit', 'View', 'Insert', 'Runtime', 'Tools', and 'Help'. A toolbar on the right contains 'Share' and 'Gemini' buttons. The left sidebar shows a file explorer with a folder named 'sample_data' containing several CSV files: 'attack_parsed_dataset.csv', 'toxicity_parsed_dataset.csv', 'twitter_parsed_dataset.csv', 'twitter_racism_parsed_dataset...', 'twitter_sexism_parsed_dataset...', and 'youtube_parsed_dataset.csv'. The main area shows the execution of three code cells:

```
[20] import pandas as pd
import torch
from torch.utils.data import Dataset
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder
from transformers import DistilBertTokenizerFast, DistilBertForSequenceClassification, Trainer, TrainingArguments
```

```
[21] !pip install -U transformers
```

```
[22] import os
```

```
[23] !pip install -U transformers datasets
os.environ["WANDB_DISABLED"] = "true"
```

Each code cell is followed by a list of requirements that are already satisfied, such as 'transformers in /usr/local/lib/python3.11/dist-packages (4.52.4)', 'filelock in /usr/local/lib/python3.11/dist-packages (from transformers) (3.18.0)', 'huggingface-hub<1.0, >=0.30.0 in /usr/local/lib/python3.11/dist-packages (from transformers) (2.0.2)', 'numpy>=1.17 in /usr/local/lib/python3.11/dist-packages (from transformers) (2.0.2)', 'packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from transformers) (24.2)', 'pyyaml>=5.1 in /usr/local/lib/python3.11/dist-packages (from transformers) (6.0.2)', 'regex!=2019.12.17 in /usr/local/lib/python3.11/dist-packages (from transformers) (2024.11.1)', 'requests in /usr/local/lib/python3.11/dist-packages (from transformers) (2.32.3)', 'tokenizers<0.22, >=0.21 in /usr/local/lib/python3.11/dist-packages (from transformers) (0.15.1)', 'safetensors>=0.4.3 in /usr/local/lib/python3.11/dist-packages (from transformers) (0.5.3)', 'tqdm>=4.27 in /usr/local/lib/python3.11/dist-packages (from transformers) (4.67.1)', 'fsspec>=2023.5.0 in /usr/local/lib/python3.11/dist-packages (from huggingface-hub<1.0, >=0.30.0) (2024.11.1)', 'typing-extensions>=3.7.4.3 in /usr/local/lib/python3.11/dist-packages (from huggingface-hub<1.0, >=0.30.0) (4.12.2)', 'hf-xet<2.0.0, >=1.1.2 in /usr/local/lib/python3.11/dist-packages (from huggingface-hub<1.0, >=0.30.0) (1.1.2)', 'charset-normalizer<4, >=2 in /usr/local/lib/python3.11/dist-packages (from requests->transformers) (3.3.2)', 'idna<4, >=2.5 in /usr/local/lib/python3.11/dist-packages (from requests->transformers) (3.10.1)', 'urllib3<3, >=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests->transformers) (2.2.3)', and 'certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests->transformers) (2024.12.14)'.

The bottom status bar shows '6:39 PM' and 'T4 (Python 3)'.

colab.research.google.com

Untitled0.ipynb - Colab

destilbert - Colab

(25) WhatsApp

Welcome To Colab - Colab

Untitled0.ipynb - Colab

Untitled0.ipynb

File Edit View Insert Runtime Tools Help

Commands

+ Code + Text

Run all

Analyze your files with code written by Gemini

Upload

logs

results

sample_data

attack_parsed_dataset.csv

toxicity_parsed_dataset.csv

twitter_parsed_dataset.csv

twitter_racism_parsed_dataset....

twitter_sexism_parsed_dataset...

youtube_parsed_dataset.csv

```

from sklearn.utils import resample
import pandas as pd

# Get the smallest class count
min_count = 1970

# Balanced list
balanced = []

# Loop through each label and resample
for label in final_df["label"].unique():
    label_df = final_df[final_df["label"] == label]
    if len(label_df) > min_count:
        resampled = resample(label_df, replace=False, n_samples=min_count, random_state=42)
    else:
        resampled = resample(label_df, replace=True, n_samples=min_count, random_state=42)
    balanced.append(resampled)

# Concatenate all balanced samples
balanced_df = pd.concat(balanced).sample(frac=1, random_state=42).reset_index(drop=True)

# Check new distribution
print("New Label Distribution:")
print(balanced_df["label"].value_counts())

```

New Label Distribution:

label

toxicity 1970

racism 1970

youtube 1970

sexism 1970

Name: count, dtype: int64

[27] # New label encoding

label_map = {'toxicity': 0, 'racism': 1, 'youtube': 2, 'sexism': 3}

Disk

72.94 GB available

Variables

Terminal

6:39 PM

T4 (Python 3)

3

The screenshot displays a Google Colab notebook titled "Untitled0.ipynb". The interface includes a top navigation bar with tabs for "Untitled0.ipynb - Colab", "destilbert - Colab", "(25) WhatsApp", "Welcome To Colab - Colab", and "Untitled0.ipynb - Colab". Below the navigation bar is a menu bar with "File", "Edit", "View", "Insert", "Runtime", "Tools", and "Help". A "Share" button and a "Gemini" icon are also present.

The left sidebar shows a "Files" panel with a folder structure:

- logs
- results
- sample_data
 - attack_parsed_dataset.csv
 - toxicity_parsed_dataset.csv
 - twitter_parsed_dataset.csv
 - twitter_racism_parsed_dataset....
 - twitter_sexism_parsed_dataset...
 - youtube_parsed_dataset.csv

The main code area contains the following Python code:

```
[27] # New label encoding
label_map = {'toxicity': 0, 'racism': 1, 'youtube': 2, 'sexism': 3}
balanced_df['label_encoded'] = balanced_df['label'].map(label_map)

import re

def clean(text):
    text = re.sub(r"http\S+", "", text)
    text = re.sub(r"[^A-Za-z\s]", "", text)
    text = re.sub(r"\s+", " ", text).strip()
    return text.lower()

balanced_df["clean_text"] = balanced_df["text"].apply(clean)

[33] # Split into train and validation sets (90/10 split)
from sklearn.model_selection import train_test_split

train_texts, val_texts, train_labels, val_labels = train_test_split(
    balanced_df['text'].tolist(),
    balanced_df['label_encoded'].tolist(),
    test_size=0.1,
    stratify=balanced_df['label_encoded'],
    random_state=42
)

# Load tokenizer
tokenizer = DistilBertTokenizerFast.from_pretrained('distilbert-base-uncased')

# Tokenize train and validation sets
train_encodings = tokenizer(train_texts, truncation=True, padding=True, max_length=128)
val_encodings = tokenizer(val_texts, truncation=True, padding=True, max_length=128)
```

The bottom status bar shows "Variables", "Terminal", "6:39 PM", and "T4 (Python 3)".

The screenshot displays a Google Colab notebook environment. The top navigation bar includes the Colab logo, the notebook name 'Untitled0.ipynb', and various utility icons like 'Share' and 'Gemini'. Below this is a menu bar with 'File', 'Edit', 'View', 'Insert', 'Runtime', 'Tools', and 'Help'. A 'Commands' bar is present with options for '+ Code', '+ Text', and 'Run all'. On the left, a 'Files' sidebar shows a directory structure with folders like 'logs', 'results', and 'sample_data', and several CSV files including 'attack_parsed_dataset.csv', 'toxicity_parsed_dataset.csv', 'twitter_parsed_dataset.csv', 'twitter_racism_parsed_dataset...', 'twitter_sexism_parsed_dataset...', and 'youtube_parsed_dataset.csv'. The main code editor area contains three code cells: [35] 'pip install -q datasets', [36] code for loading datasets from dictionaries and creating Dataset objects, and [37] code for tokenizing the datasets and renaming the label column. The bottom status bar shows '6:39 PM' and 'T4 (Python 3)'.

colab.research.google.com

Untitled0.ipynb - Colab

destilbert - Colab

(25) WhatsApp

Welcome To Colab - Colab

Untitled0.ipynb - Colab

File Edit View Insert Runtime Tools Help

Commands + Code + Text Run all

Files

Analyze your files with code written by Gemini Upload

logs

results

sample_data

attack_parsed_dataset.csv

toxicity_parsed_dataset.csv

twitter_parsed_dataset.csv

twitter_racism_parsed_dataset...

twitter_sexism_parsed_dataset...

youtube_parsed_dataset.csv

[35] !pip install -q datasets

[36] from datasets import Dataset

```
# Rebuild train and validation into HF-compatible Dataset format
train_dict = {
    'text': train_texts,
    'label': train_labels
}

val_dict = {
    'text': val_texts,
    'label': val_labels
}

train_dataset = Dataset.from_dict(train_dict)
val_dataset = Dataset.from_dict(val_dict)
```

[37] def tokenize_fn(example):

```
    return tokenizer(example['text'], truncation=True, padding='max_length', max_length=128)

# Tokenize datasets
train_dataset = train_dataset.map(tokenize_fn, batched=True)
val_dataset = val_dataset.map(tokenize_fn, batched=True)

# Rename label column for Trainer compatibility
train_dataset = train_dataset.rename_column("label", "labels")
val_dataset = val_dataset.rename_column("label", "labels")

# Set format to PyTorch
train_dataset.set_format(type='torch', columns=['input_ids', 'attention_mask', 'labels'])
val_dataset.set_format(type='torch', columns=['input_ids', 'attention_mask', 'labels'])
```

Disk 72.94 GB available

Variables Terminal

6:39 PM T4 (Python 3)

colab.research.google.com

Untitled0.ipynb - Colab

File Edit View Insert Runtime Tools Help

Commands + Code + Text Run all

Files

- ..
- logs
- results
- sample_data
 - attack_parsed_dataset.csv
 - toxicity_parsed_dataset.csv
 - twitter_parsed_dataset.csv
 - twitter_racism_parsed_dataset...
 - twitter_sexism_parsed_dataset...
 - youtube_parsed_dataset.csv

Disk 72.94 GB available

Map: 100% 7092/7092 [00:07<00:00, 1001.17 examples/s]

Map: 100% 788/788 [00:00<00:00, 1085.10 examples/s]

```
[38] from transformers import DistilBertForSequenceClassification# Number of unique labels (classes)num_labels = len(set(tr
```

```
from transformers import TrainingArguments
training_args = TrainingArguments(
    output_dir='./results',
    num_train_epochs=4,
    per_device_train_batch_size=16,
    per_device_eval_batch_size=32,
    warmup_steps=100,
    weight_decay=0.01,
    logging_dir='./logs',
    logging_steps=10,
    save_strategy="no"
)
```

Using the 'WANDB_DISABLED' environment variable is deprecated and will be removed in v5. Use the --report_to flag to con

```
[41] from sklearn.metrics import accuracy_score, f1_score

def compute_metrics(p):
    preds = p.predictions.argmax(axis=1)
    labels = p.label_ids
    return {
        "accuracy": accuracy_score(labels, preds),
        "macro_f1": f1_score(labels, preds, average='macro')
    }
```

Variables Terminal

6:39 PM T4 (Python 3)

colab.research.google.com

Untitled0.ipynb - Colab

File Edit View Insert Runtime Tools Help

Commands + Code + Text Run all

Files

Analyze your files with code written by Gemini Upload

logs

results

sample_data

attack_parsed_dataset.csv

toxicity_parsed_dataset.csv

twitter_parsed_dataset.csv

twitter_racism_parsed_dataset....

twitter_sexism_parsed_dataset...

youtube_parsed_dataset.csv

Disk 72.94 GB available

```
[41]
trainer = Trainer(
    model=model,
    args=training_args,
    train_dataset=train_dataset,
    eval_dataset=val_dataset,
    compute_metrics=compute_metrics
)
```

trainer.train()

1480	0.203000
1490	0.238300
1500	0.294900
1510	0.265400
1520	0.226600
1530	0.187200
1540	0.231000
1550	0.248700
1560	0.181200
1570	0.314200
1580	0.246600
1590	0.193900
1600	0.237500
1610	0.344900
1620	0.238900

Variables Terminal

6:39 PM T4 (Python 3)

colab.research.google.com

Untitled0.ipynb - Colab

File Edit View Insert Runtime Tools Help

Commands + Code + Text Run all

Files

- ..
- logs
- results
- sample_data
- attack_parsed_dataset.csv
- toxicity_parsed_dataset.csv
- twitter_racism_parsed_dataset...
- twitter_sexism_parsed_dataset...
- youtube_parsed_dataset.csv

1770 0.215400

```
TrainOutput(global_step=1776, training_loss=0.3667051704885723, metrics={'train_runtime': 303.3996,
'train_samples_per_second': 93.5, 'train_steps_per_second': 5.854, 'total_flos': 939492299096064.0, 'train_loss':
0.3667051704885723, 'epoch': 4.0})
```

[45] # Evaluate model on validation set
trainer.evaluate()

[25/25 00:02]

```
{'eval_loss': 0.6470827460289001,
'eval_accuracy': 0.7373096446700508,
'eval_macro_f1': 0.7378934462747743,
'eval_runtime': 2.8971,
'eval_samples_per_second': 271.999,
'eval_steps_per_second': 8.629,
'epoch': 4.0}
```

from transformers import DistilBertTokenizer
tokenizer = DistilBertTokenizer.from_pretrained('distilbert-base-uncased')

```
val_encodings = tokenizer(
    val_texts,
    padding="max_length",
    truncation=True,
    max_length=128,
    return_tensors="pt"
```

[48] import torch

```
val_labels_tensor = torch.tensor(val_labels)
```

[50] from torch.utils.data import Dataset

Disk 72.94 GB available

Variables Terminal

6:39 PM T4 (Python 3)

colab.research.google.com

Untitled0.ipynb - Colab

File Edit View Insert Runtime Tools Help

Commands + Code + Text Run all

Files

- ..
- logs
- results
- sample_data
 - attack_parsed_dataset.csv
 - toxicity_parsed_dataset.csv
 - twitter_parsed_dataset.csv
 - twitter_racism_parsed_dataset....
 - twitter_sexism_parsed_dataset...
 - youtube_parsed_dataset.csv

Analyze your files with code written by Gemini Upload

```
[50] from torch.utils.data import Dataset

class ValDataset(Dataset):
    def __init__(self, encodings, labels):
        self.encodings = encodings
        self.labels = labels

    def __len__(self):
        return len(self.labels)

    def __getitem__(self, idx):
        return {
            'input_ids': self.encodings['input_ids'][idx],
            'attention_mask': self.encodings['attention_mask'][idx],
            'label': self.labels[idx]
        }

[51] from torch.utils.data import DataLoader

val_loader = DataLoader(val_dataset, batch_size=32)

[53] import torch

device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
print("Using device:", device)

Using device: cuda

[56] all_preds = []
all_labels = []
```

Disk 72.94 GB available

Variables Terminal

6:39 PM T4 (Python 3)

colab.research.google.com

Untitled0.ipynb - Colab

File Edit View Insert Runtime Tools Help

Commands + Code + Text Run all

Files

- ..
- logs
- results
- sample_data
- attack_parsed_dataset.csv
- toxicity_parsed_dataset.csv
- twitter_parsed_dataset.csv
- twitter_racism_parsed_dataset...
- twitter_sexism_parsed_dataset...
- youtube_parsed_dataset.csv

```
[56] all_preds = []
      all_labels = []

      model.to(device)
      model.eval()

      with torch.no_grad():
          for batch in val_loader:
              input_ids = batch['input_ids'].to(device)
              attention_mask = batch['attention_mask'].to(device)

              label_key = 'label' if 'label' in batch else 'labels'
              labels = batch[label_key].to(device)

              outputs = model(input_ids=input_ids, attention_mask=attention_mask)
              logits = outputs.logits
              preds = torch.argmax(logits, dim=1)

              all_preds.extend(preds.cpu().numpy())
              all_labels.extend(labels.cpu().numpy())

[57] from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
      import matplotlib.pyplot as plt

      # Define class names (edit if needed)
      label_names = ['toxicity', 'racism', 'youtube', 'sexism']

      # Generate the confusion matrix
      cm = confusion_matrix(all_labels, all_preds)
      disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=label_names)

      # Plot
      plt.figure(figsize=(8, 6))
      disp.plot(cmap='Blues', values_format='d')
      plt.title("Confusion Matrix")
```

Disk 72.94 GB available

Variables Terminal

6:39 PM T4 (Python 3)



