

МИНИСТЕРСТВО ЦИФРОВОГО РАЗВИТИЯ, СВЯЗИ И МАССОВЫХ
КОММУНИКАЦИЙ РОССИЙСКОЙ ФЕДЕРАЦИИ
Ордена Трудового Красного Знамени федеральное государственное
бюджетное образовательное учреждение высшего образования
**«Московский Технический Университет Связи И Информатики
(MTUCI)»**

Кафедра «Математическая кибернетика и информационные технологии»

Лабораторная Работа 3

по дисциплине

«Машинное обучение»

Выполнил: студент 3 курса гр. БВТ2201
Ньяти Каелиле

Москва 2025 г

Содержание

1. Задание
2. Ход работы
3. Вывод

1. Задание

Написать RNN/GRU/LSTM (любую из) модель которая будет способна классифицировать текст на токсичный или не токсичный. Датасет можно выбрать любой, но если лень искать то можете взять этот. Аккуратно выбирать размеры модели так как у RNN моделей большие проблемы со скоростью обучения. В качестве токенизатора можно взять любой токенизатор с hf либо написать свой эмбединг слой который будет превращать именно слова в эмбединги, тут по желанию.

В качестве лосса это BCE. В качестве оптимизатор Adam/AdamW, lr=3e-4. Но если захотите что другое то можно другое. Loss можете взвесить так как явный дисбаланс классов. По метрикам считайте ассурасу, f1, roc-auc, pr-auc. Не забудьте проводить расчеты на cuda.

2. Ход работы

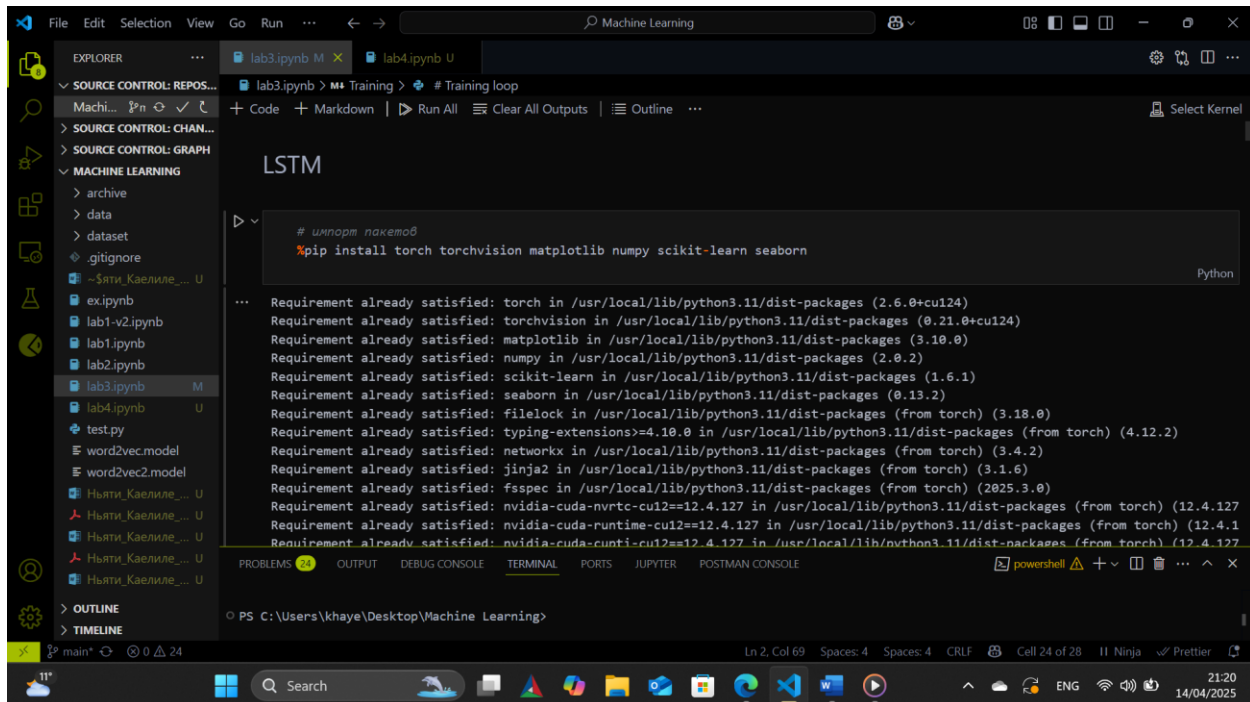


Рис 1. Библиотеки

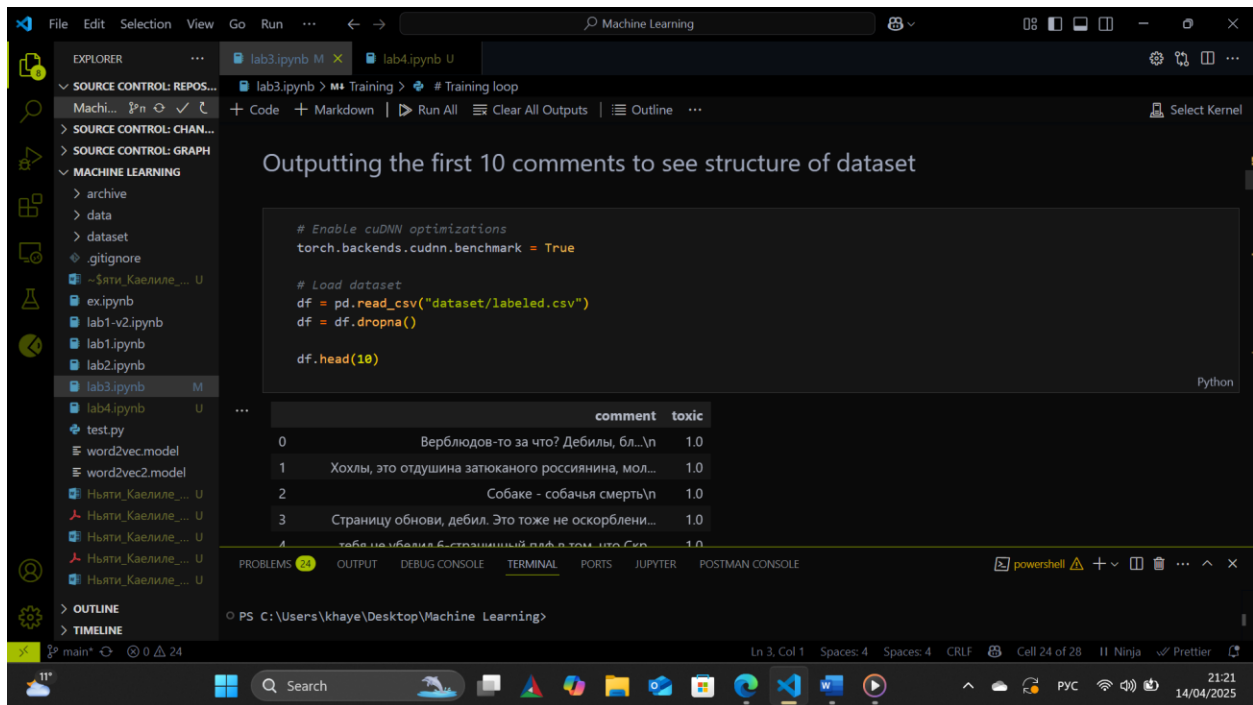
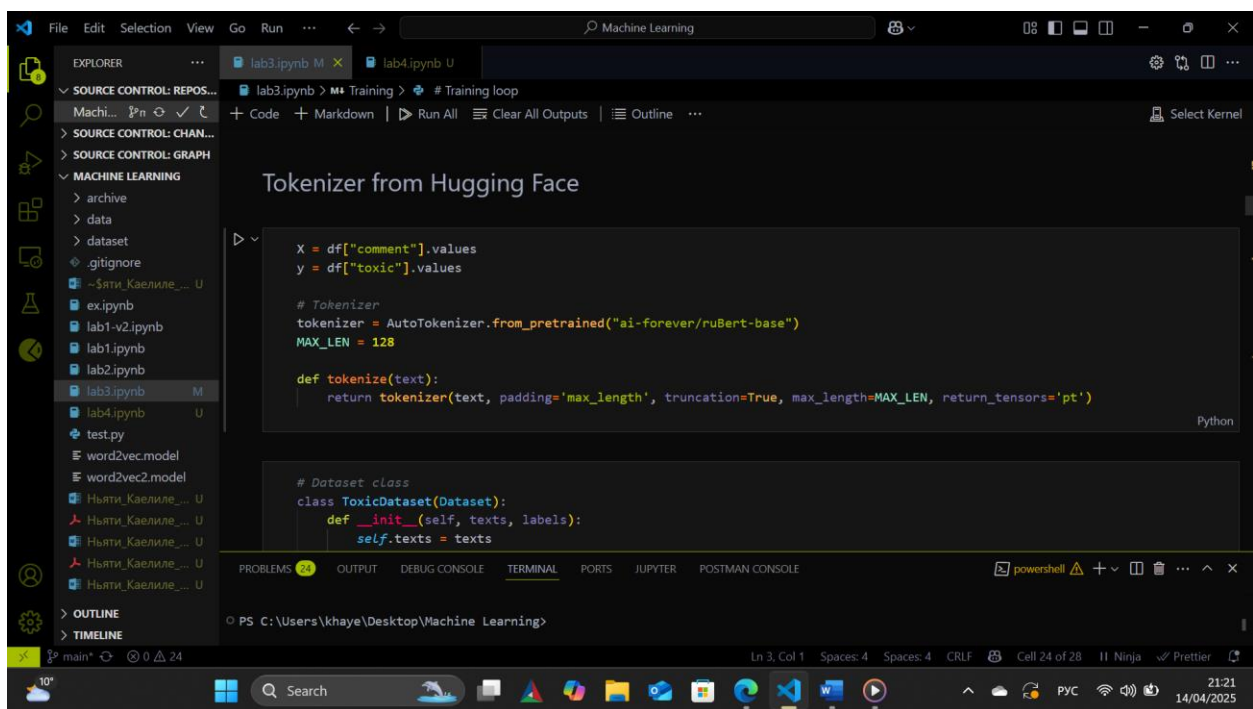


Рис 2. Дата



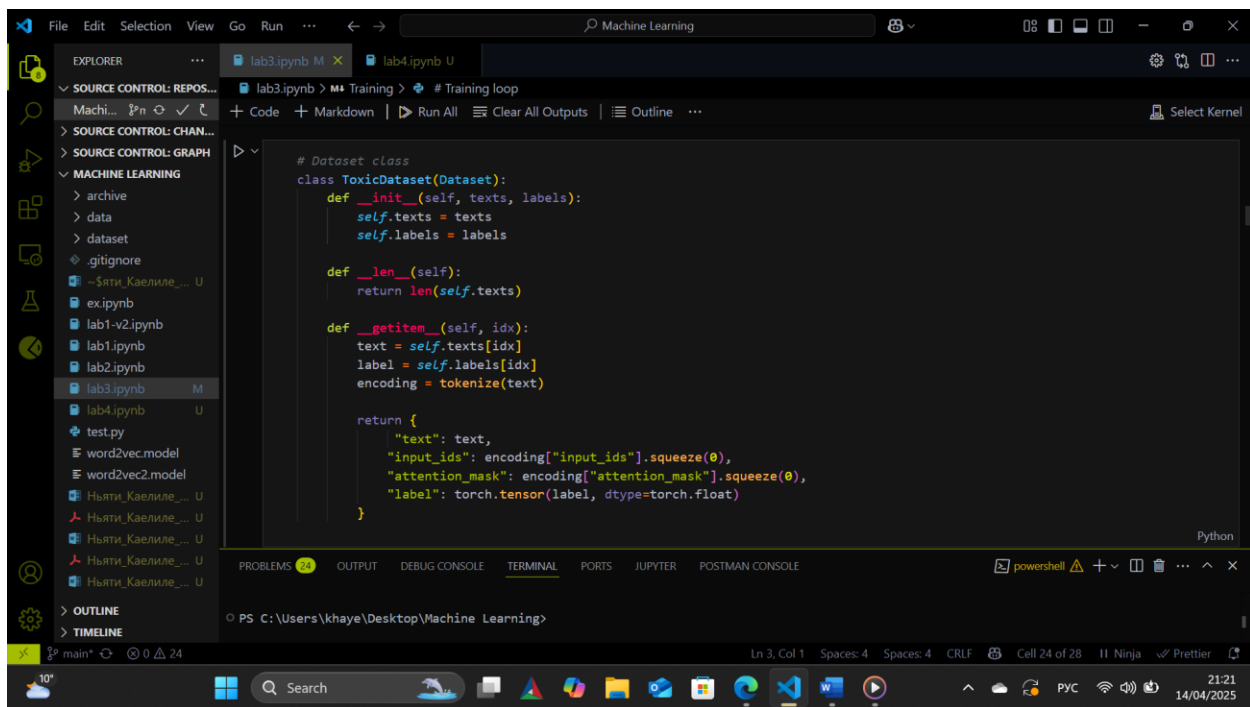
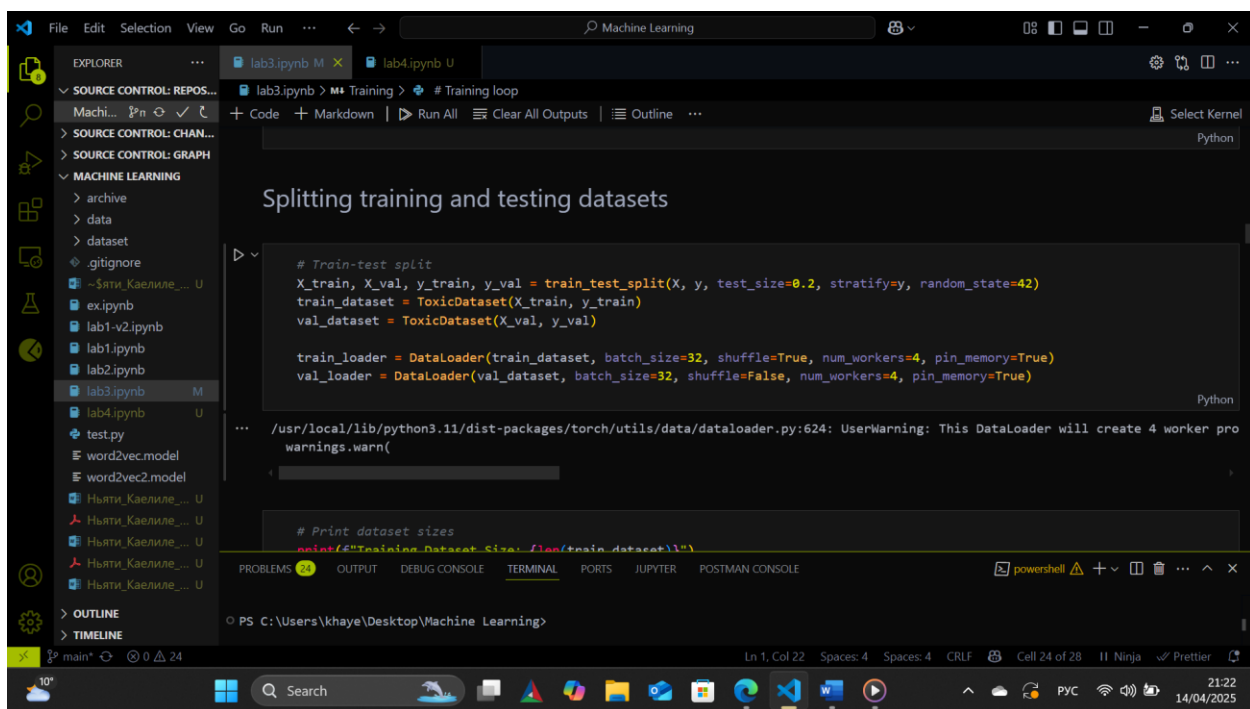
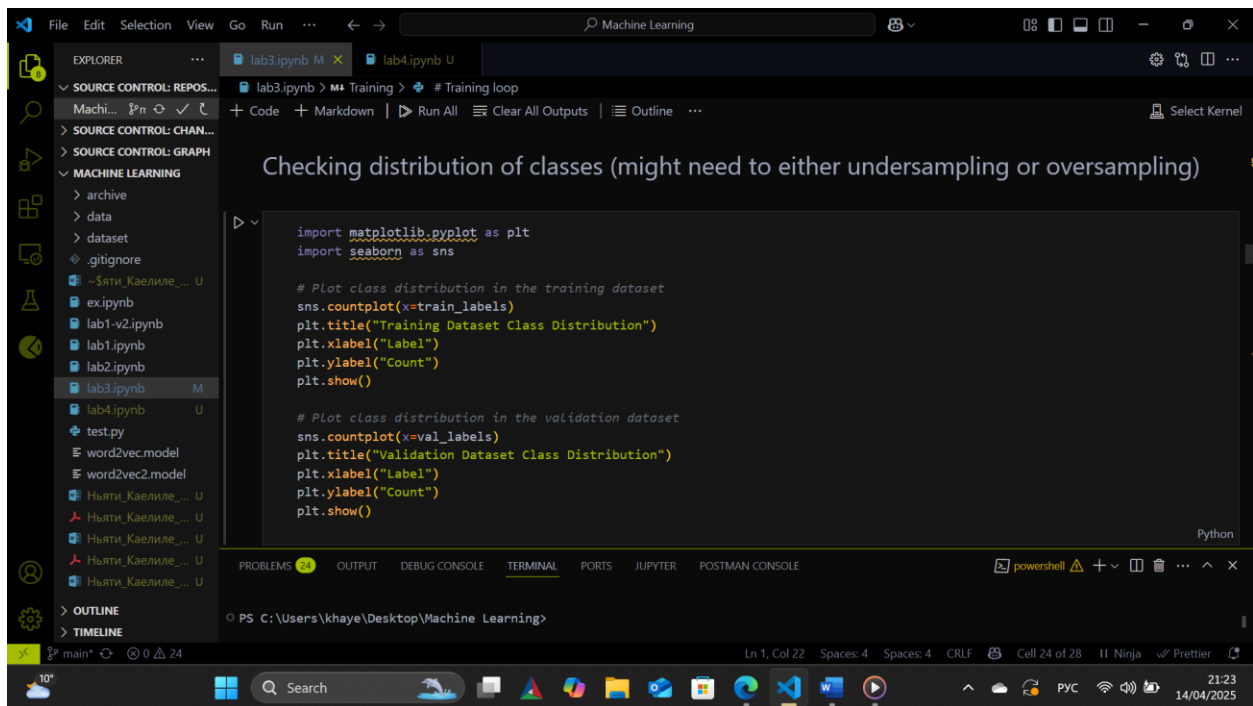
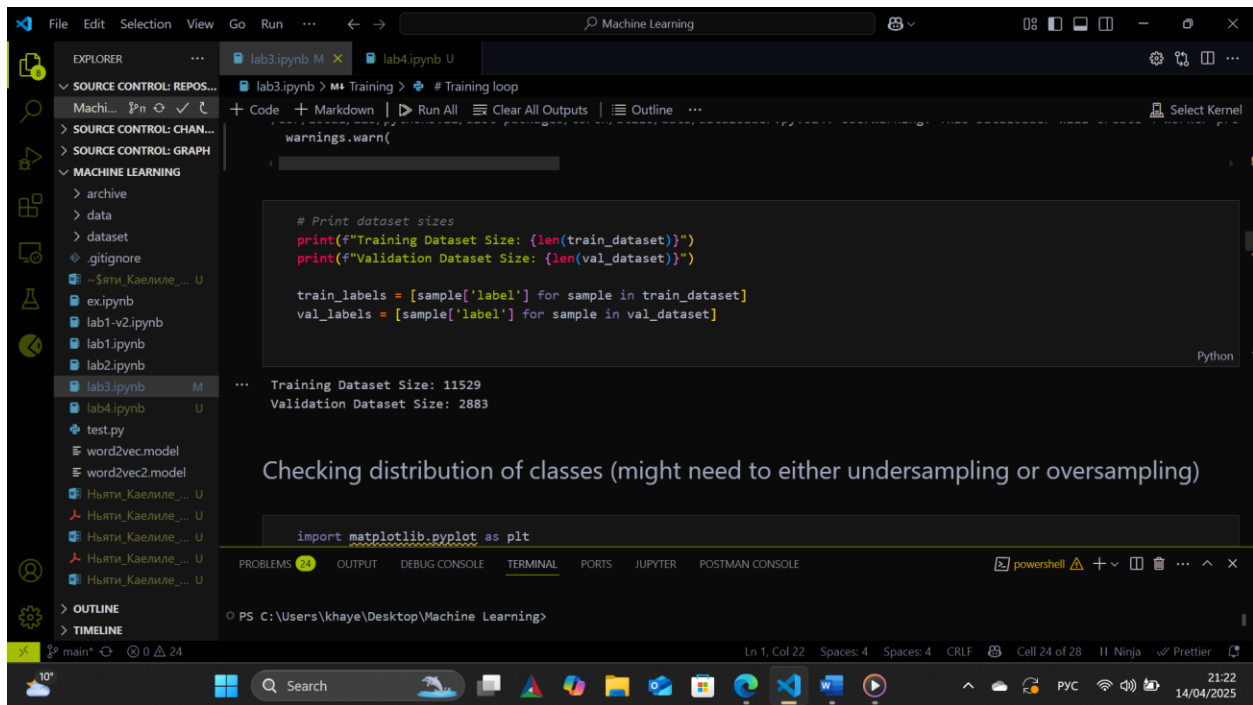
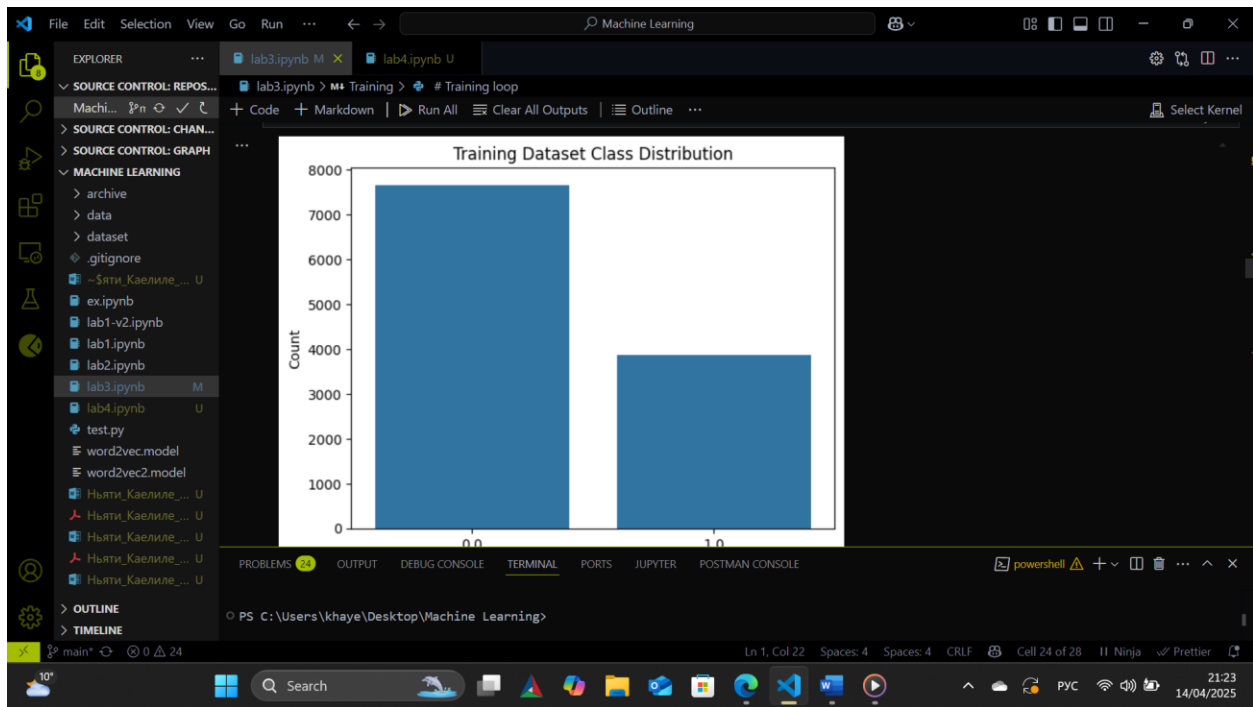
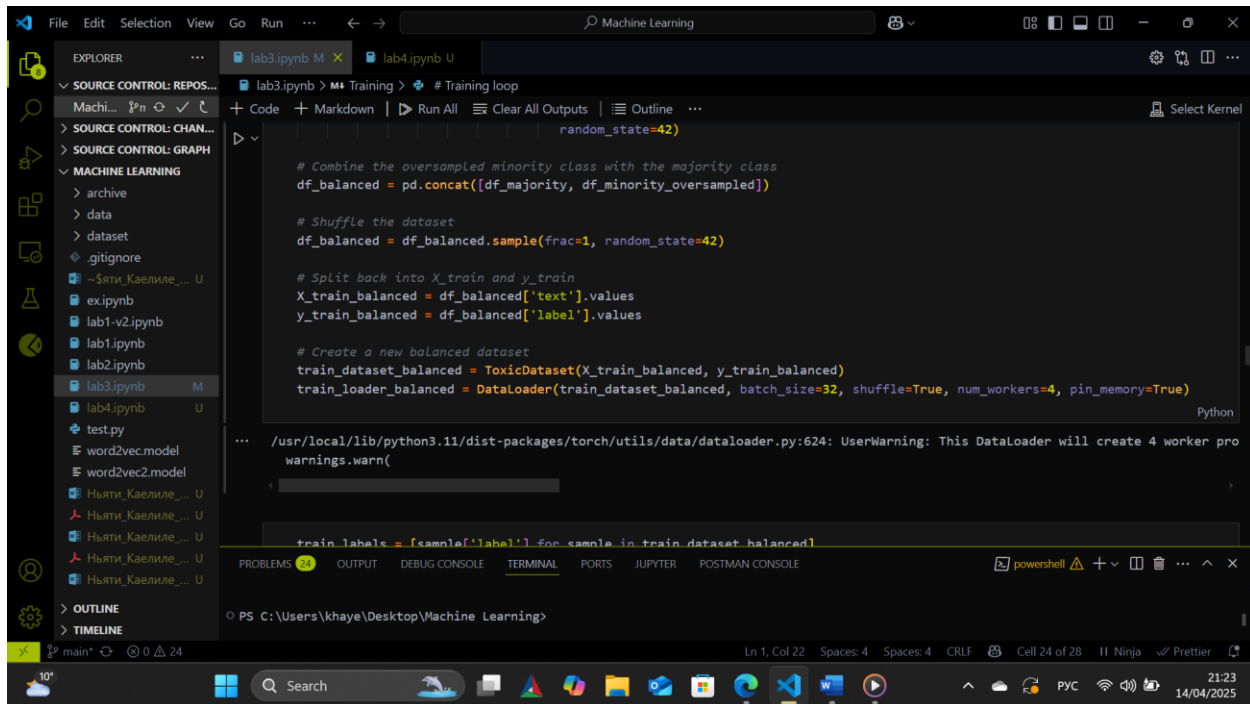
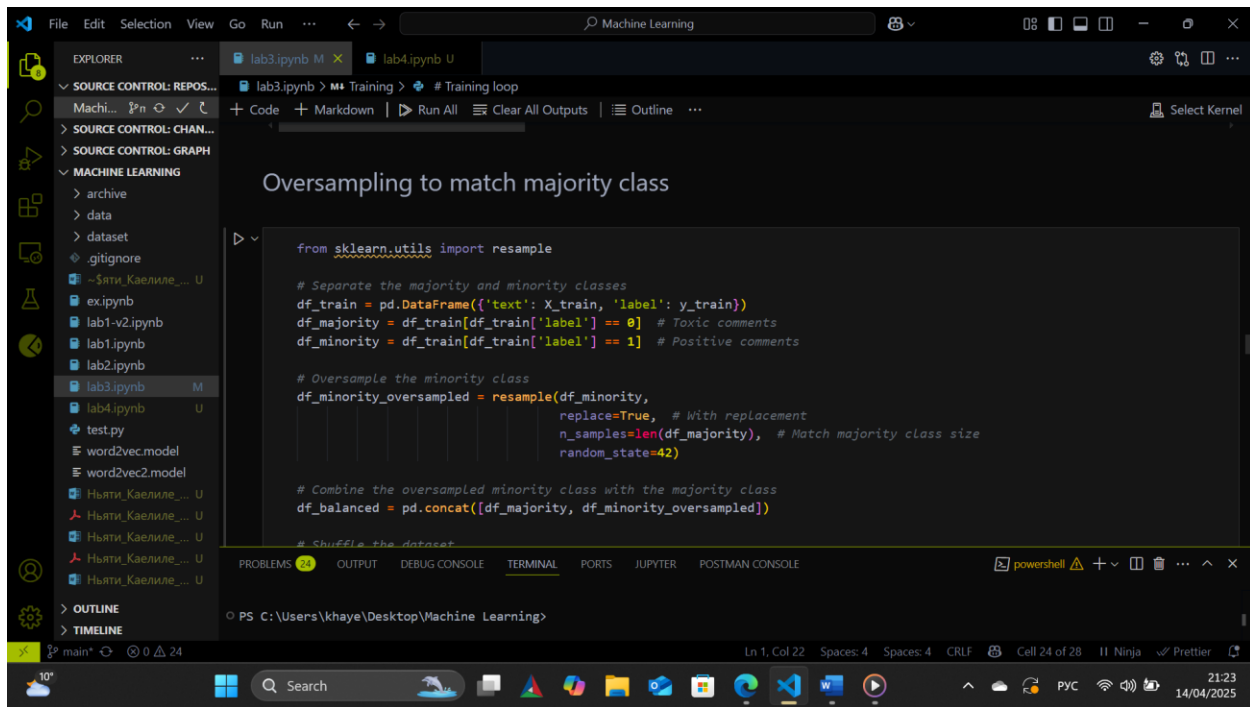


Рис 3. Токенизатор









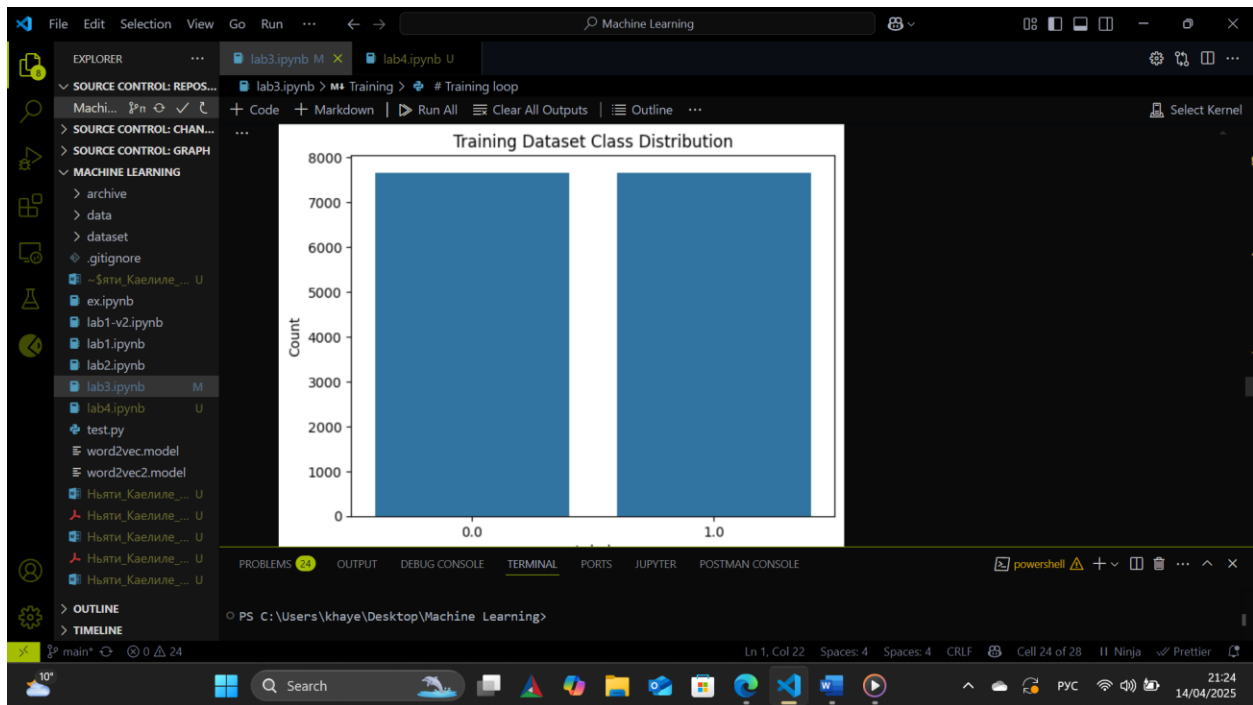


Рис 4. Oversampling

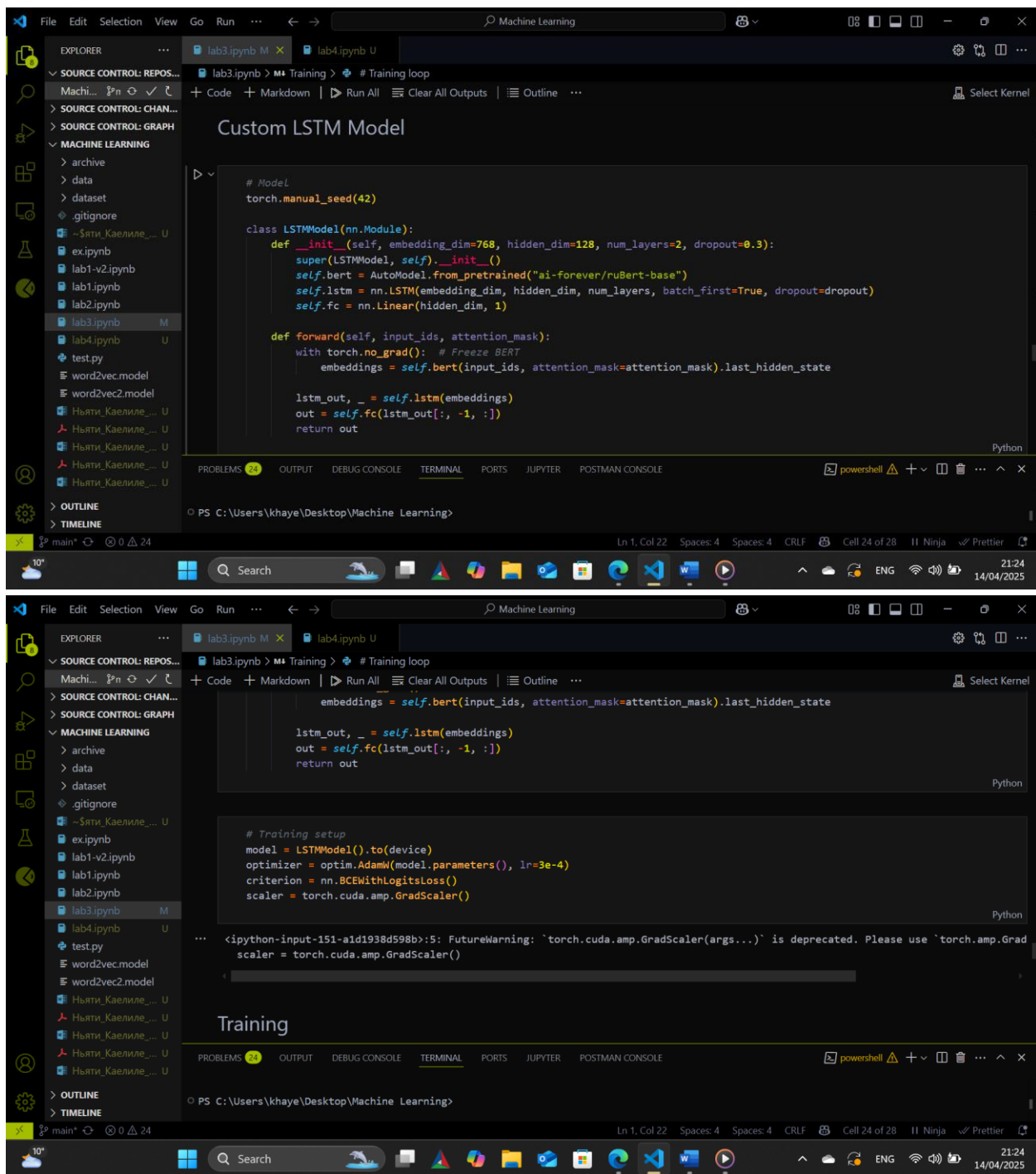


Рис 5. LSTM model

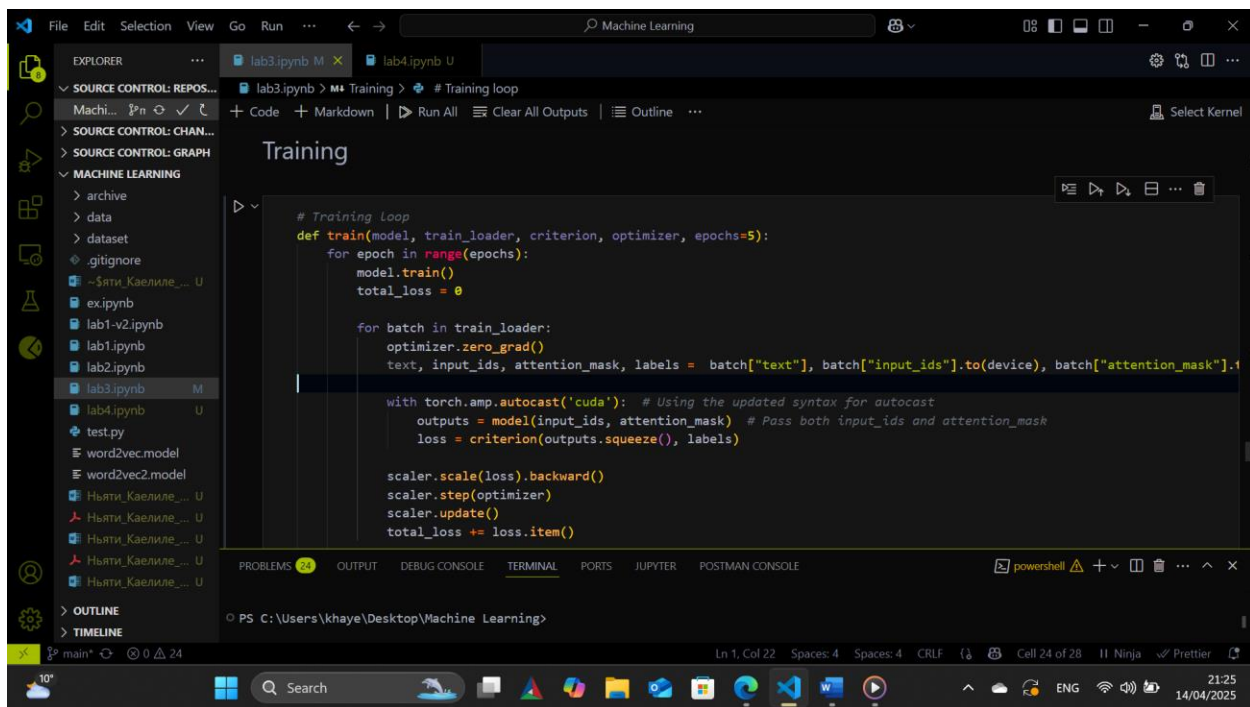
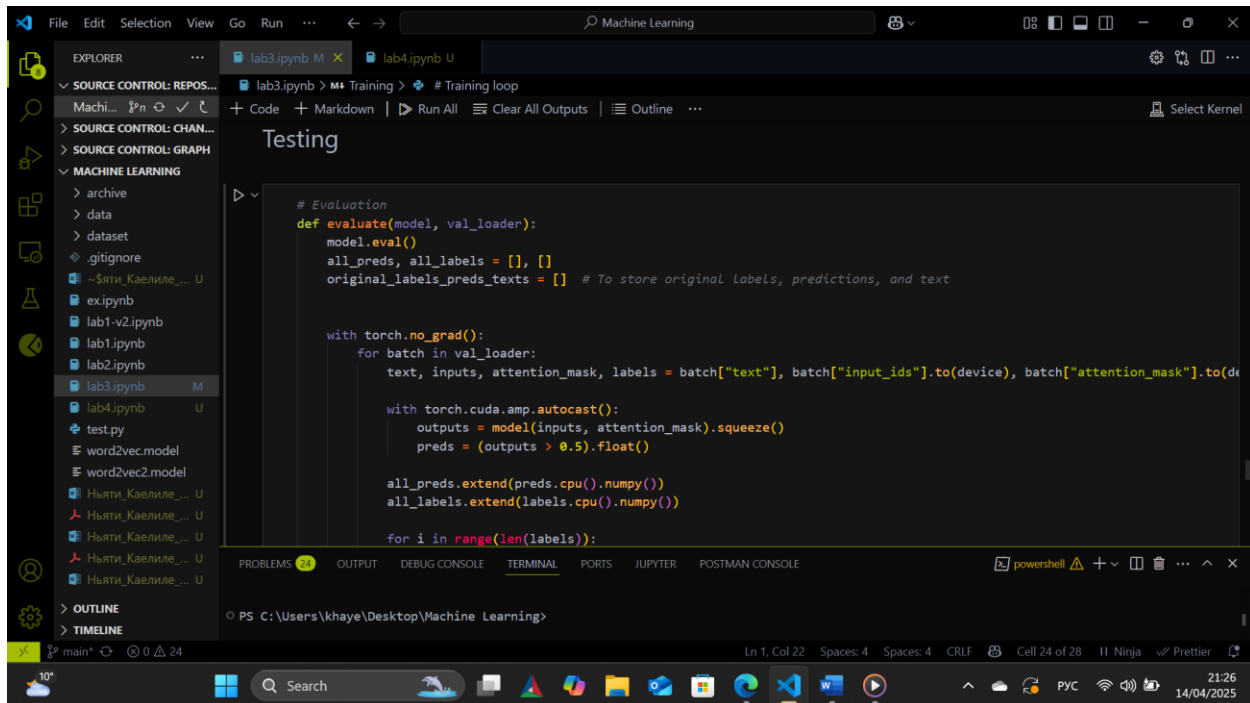


Рис 6. Тренирование



The screenshot shows the VS Code editor with a Jupyter Notebook open. The Explorer sidebar on the left shows a project structure with files like `lab3.ipynb`, `lab4.ipynb`, `test.py`, and `word2vec.model`. The active cell in the notebook contains Python code for a training loop. The code includes data batching, model inference with `torch.cuda.amp.autocast()`, and calculation of accuracy, F1 score, ROC-AUC, and PR-AUC. It also prints the top 10 original labels, predictions, and text.

```
# Training loop
text, inputs, attention_mask, labels = batch["text"], batch["input_ids"].to(device), batch["attention_mask"].to(device)

with torch.cuda.amp.autocast():
    outputs = model(inputs, attention_mask).squeeze()
    preds = (outputs > 0.5).float()

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

for i in range(len(labels)):
    original_labels_preds_texts.append((labels[i].item(), preds[i].item(), text[i]))

acc = accuracy_score(all_labels, all_preds)
f1 = f1_score(all_labels, all_preds)
roc_auc = roc_auc_score(all_labels, all_preds)
pr_auc = average_precision_score(all_labels, all_preds)

print(f"Accuracy: {acc}, F1: {f1}, ROC-AUC: {roc_auc}, PR-AUC: {pr_auc}")

# Print top 10 original Labels, predictions, and text
print("\nTop 10 Original Labels, Predictions, and Text:")
for i, (label, pred, text) in enumerate(original_labels_preds_texts[:10]):
```

The screenshot shows the VS Code editor with a Jupyter Notebook open. The Explorer sidebar on the left shows a project structure with files like `lab3.ipynb`, `lab4.ipynb`, `test.py`, and `word2vec.model`. The active cell in the notebook contains Python code for training and evaluation. The code includes a `train` function call and a `evaluate` function call. The output of the training process is visible in the terminal, showing the loss for each of the 5 epochs.

```
print(f"sample {i+1}: ")
print(f" Text: {text}")
print(f" Label: {label}, Prediction: {pred}")
print("-" * 50) # Separator for readability

train(model, train_loader_balanced, criterion, optimizer)

...

/usr/local/lib/python3.11/dist-packages/torch/utils/data/dataloader.py:624: UserWarning: This DataLoader will create 4 worker processes
warnings.warn(
Epoch 1, Loss: 0.2835847528108085
Epoch 2, Loss: 0.21760987889332076
Epoch 3, Loss: 0.19031515084789136
Epoch 4, Loss: 0.16063627262677377
Epoch 5, Loss: 0.14183037186739966

evaluate(model, val_loader)
```

