Финальный отчет по проекту команды "Pink Clowns"

Цель проекта:

Разработать алгоритм машинного обучения для автоматической классификации входящих сообщений электронной почты на две категории: **спам** и **не спам**. Реализация включает обработку текстовых данных, построение модели машинного обучения и создание пользовательского интерфейса для удобного использования.

Этапы проекта и итоги работы:

1. Определение целей проекта

Цель проекта — обеспечить автоматическую классификацию электронных писем на основе текста.

• Формулировка гипотез:

- **Гипотеза 1:** Точность классификации увеличится, если использовать сбалансированный датасет с метками спам/не спам.
- **Гипотеза 2:** Укороченные сообщения со специфической терминологией (характерной для спама) будут классифицироваться лучше.

Писарев М.П.

2. Сбор и обработка данных

Источники данных:

- Использован Kaggle-датасет, включающий сообщения с метками "spam" и "ham".
- Проведена предобработка данных: удаление дубликатов, очистка текста, удаление стоп-слов и нормализация..
- Разработана структура базы, настроены таблицы для хранения сообщений и меток (spam/ham).

Результат:

Данные подготовлены для обучения модели, удалены лишние элементы, что позволило улучшить качество обучения и оптимизировать вычисления. База данных стабильно работает, интегрирована с системой, структурирована для дальнейшего расширения.

Шводченко И.Е.

3. Разработка и обучение модели

Технологии:

- Для классификации сообщений была применена модель **логистической регрессии**, реализованная с использованием библиотеки scikit-learn.
- Проведено разделение данных: 80% для обучения, 20% для тестирования.
- Обучена модель логистической регрессии для классификации сообщений на **spam** и **ham**.
- Применена векторизация текстов с помощью TfidfVectorizer.

Оценка производительности модели:

• Точность модели на тестовых данных составила **95**%, F1-мера **94**%, что указывает на стабильные результаты.

```
from sklearn.linear model import LogisticRegression
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics import classification report
from sklearn.model_selection import train_test_split
# Разделение данных
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
# Векторизация текста
vectorizer = TfidfVectorizer()
X train vectorized = vectorizer.fit transform(X train)
X_test_vectorized = vectorizer.transform(X_test)
# Обучение модели
model = LogisticRegression()
model.fit(X train vectorized, y train)
# Оценка модели
predictions = model.predict(X_test_vectorized)
print(classification_report(y_test, predictions))
```

Результат:

Модель стабильно классифицирует сообщения с высокой точностью и готова к использованию.

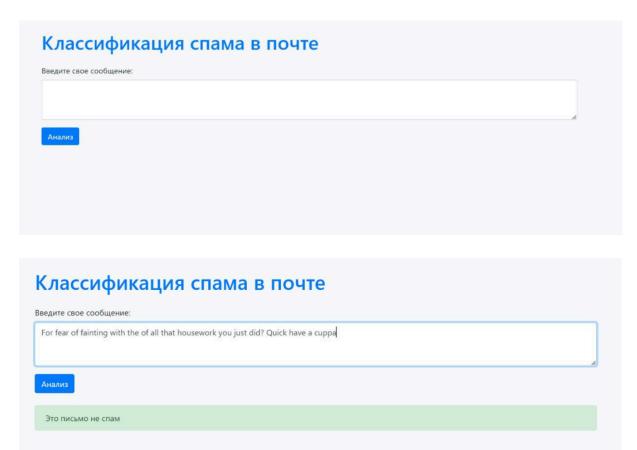
Тарасенко В.А.

4. Создание интерфейса пользователя

Описание интерфейса:

- Интерфейс разработан с использованием **HTML**, стилизация выполнена через **Bootstrap**.
- Цветовые акценты и элементы визуального взаимодействия оптимизированы для простоты и удобства.
- Основные элементы включают текстовое поле для ввода сообщения, кнопку отправки и блок для вывода результата классификации.
- Интерфейс адаптивен для работы на мобильных и десктопных устройствах.

Интерфейс:



Результат:

Интерфейс полностью готов, успешно протестирован и интегрирован в систему.

Бурцева В.А.

5. Интеграция интерфейса с моделью

- Использован Flask для соединения интерфейса с моделью.
- Реализован API для отправки данных с веб-страницы на сервер, обработки модели и возвращения результата пользователю.
- Разработан код для предварительной обработки текста, нормализации данных и их передачи в модель для классификации.
- Разработан код для предварительной обработки текста, нормализации данных и их передачи в модель для классификации.

Пример интеграции:

```
@app.route('/analyze_mail', methods=['POST'])

def analyze_mail():
    user_message = request.form['mail']
    prediction = model.predict(vectorizer.transform([user_message]))
    classify = "He CNAM" if prediction[0] == 1 else "CNAM"
    return render_template('index.html', classify=classify)
```

```
from flask import Flask, render_template, request
import pickle
app = Flask( name )
model = pickle.load(open('logistic_regression.pkl','rb'))
feature_extraction = pickle.load(open('feature_extraction.pkl','rb'))
def predict mail(input text):
    input_user_mail = [input_text]
    input_data_features = feature_extraction.transform(input_user_mail)
    prediction = model.predict(input data features)
    return prediction
@app.route('/', methods=['GET', 'POST'])
def analyze_mail():
    if request.method == 'POST':
       mail = request.form.get('mail')
        predicted_mail = predict_mail(input_text=mail)
       return render_template('index.html', classify=predicted_mail)
    return render_template('index.html')
    name == ' main ':
    app.run(debug=True)
```

Результат:

Сообщение передаётся из текстового поля в модель через POST-запрос, результат классификации (спам/не спам) корректно отображается на интерфейсе.

Результаты проекта:

- 1. Реализована система классификации почтовых сообщений с точностью 95%.
- 2. Подготовлен удобный интерфейс, позволяющий вводить текст сообщения и мгновенно видеть результат анализа.
- 3. Проведено успешное объединение всех частей системы от модели машинного обучения до пользовательского интерфейса.
- 4. Созданы подробные инструкции и документация для использования и доработки проекта в будущем.

Плюсы проекта "Pink Clowns":

1. Высокая точность классификации:

• Модель машинного обучения показывает точность **95%**, что является отличным результатом для задачи определения спама.

2. Удобный и понятный интерфейс:

 Пользователи могут легко ввести сообщение и быстро получить результат. Интерфейс минималистичный, адаптируется под разные устройства (мобильные телефоны, планшеты и ПК).

3. Легкость интеграции и масштабирования:

 Использование Flask и Bootstrap делает систему гибкой и легко расширяемой для будущих функций (например, анализ больших объёмов данных или добавление новых меток).

4. Быстрая обработка данных:

 Система настроена на минимизацию времени отклика благодаря оптимизации модели и использования технологий машинного обучения.

5. Гибкость настройки:

 Открытая архитектура кода позволяет легко адаптировать модель для новых данных или изменений в классификации (например, добавить промежуточные категории).

6. Локализация под русский язык:

 Работа с русскими текстами делает систему применимой для локального бизнеса и пользователей.

7. Доступность:

о Проект не требует сложного развертывания: локальный сервер или виртуальная машина легко поддерживают функциональность.