

УДК 004.89

МЕТОД АУГМЕНТАЦИИ ТЕКСТОВЫХ ДАННЫХ С СОХРАНЕНИЕМ ЭМОЦИОНАЛЬНОЙ ОКРАСКИ

Матвеева А.А. (ИТМО),

Научный руководитель – кандидат технических наук, доцент, Махныткина О.В. (ИТМО)

Введение. Исследования выполнены за счет финансирования университета ИТМО в рамках НИР №623088 «Разработка русскоязычного персонифицированного эмоционального диалогового агента».

Эмоциональные диалоговые агенты представляют собой инновационный инструмент для улучшения качества человеко-машинного взаимодействия, они способны поддерживать диалог на эмоциональном уровне, откликаться на настроение пользователя и показывать сочувствие. Одним из ограничений при их разработке является отсутствие больших наборов данных, содержащих эмоциональные диалоговые данные, особенно на русском языке. Одним из решений данной проблемы может стать аугментация текстовых данных. В данной работе предложен метод аугментации текстовых данных на основе GPT [2] (GPT-3.5-turbo, RUGPT3, YandexGPT) моделей, позволяющий сохранять эмоциональную окраску исходной фразы.

Основная часть. Аугментация текстовых данных с возможностью сохранения эмоциональной окраски исходной фразы проводилась с использованием GPT моделей: GPT-3.5-turbo; RUGPT3[1]; YandexGPT. На первом этапе были выделены способы передачи эмоционального состояния в письменной речи. Были определены ключевые слова, фразы, сочетания знаков пунктуации, которые отражают эмоциональное состояние. На втором этапе была разработана методика формирования промптов для GPT моделей с использованием выявленных способов передачи эмоций в письменной речи. Оценка сохранения эмоциональной окраски реплик, сгенерированных рассматриваемыми моделями, осуществлялась с помощью классификатора эмоций на основе модели BERT. Эксперименты по аугментации данных с возможностью сохранения эмоциональной окраски проводились для текстов на русском и английских языках. Для данного исследования были выбраны два свободно распространяемых датасета MELD [3] для английского языка и RESD (Russian Emotional Speech Dialogues) для русского.

«MELD» - набор данных на английском языке, содержит более 1400 диалогов и 13000 высказываний из сериала «Друзья». Каждое высказывание в диалоге было аннотировано одной из семи эмоций: гнев, отвращение, печаль, радость, нейтральность, удивление и страх.

«RESD» - набор данных на русском языке, содержит 1396 высказываний, каждое из которых аннотировано одной из семи эмоций: гнев, отвращение, страх, энтузиазм, счастье, нейтральность и печаль.

Выводы. В данной работе был проведен сравнительный анализ использования различных GPT моделей для задачи аугментации текстовых данных с возможностью сохранения эмоциональной окраски.

Список использованных источников:

1. Konodyuk N., Tikhonova M. Continuous Prompt Tuning for Russian: How to Learn Prompts Efficiently with RuGPT3? // In: Burnaev, E., et al. Recent Trends in Analysis of Images, Social Networks and Texts. AIST 2021. Communications in Computer and Information Science. – 2021. – №1573. – P. 30–40.
2. Floridi L., Chiriatti M. GPT-3: Its Nature, Scope, Limits, and Consequences // Minds & Machines. – 2020. – №30. – P. 681–694.
3. Chen S.Y., Hsu C.C., Kuo C.C., Ku L.W. EmotionLines: An Emotion Corpus of Multi-Party Conversations // In Proceedings of the Eleventh International Conference on Language Resources and Evaluation – 2018. – №1. – P. 1597-1601.