

题目：毒性评论严重程度评估

成员：

杨汉伟 2006500035

龙 豪 2006500003

陈弘耀 2006500007

Kaggle 链接：<https://www.kaggle.com/competitions/jigsaw-toxic-severity-rating/overview>

科学意义：

随着互联网的发展，社交媒体平台上的评论数量呈现爆发式增长，这也导致了大量的网络暴力、仇恨言论、性别歧视等不良现象的出现。如何评估和过滤这些评论成为了亟待解决的问题。而毒性评论的严重程度评估正是解决这一问题的重要环节。该问题的解决有助于社交媒体平台建立更加友好、安全的社区氛围，遏制不良言论的传播。

科学问题：

本次比赛旨在利用机器学习技术，根据给定的文本评论，判断该评论的毒性程度。具体来说，需要预测该评论是否含有以下几种毒性因素：恶意、严重的恶意、仇恨、不敬、性或性别暴力、威胁。该问题的解决有助于社交媒体平台快速、准确地过滤掉毒性评论，保障用户的安全和体验。

研究内容：

本次比赛提供的数据集包含大量的英文文本评论，同时还提供了每个评论对应的6种毒性因素的标注。我们计划使用自然语言处理和深度学习算法对该数据集进行建模和训练，并对评论的毒性进行分类。具体来说，我们计划使用以下研究方法：

语言：python

模型：文本特征提取：词袋模型、TF-IDF 模型、词嵌入模型等 训练和优化：

使用 Scikit-learn、Keras、TensorFlow 等机器学习库，训练和优化 多个不同的分类器模型，包括逻辑回归、决策树、支持向量机、深度神经网络等

预期目标：该项目旨在利用自然语言处理和机器学习算法构建一个文本分类器，能够自动检测并评估给定评论的毒性等级。该分类器应该能够对新的评论进行准确的分类，同时保证在处理大量数据时具有高效性和稳定性；

附数据介绍：

数据集：该数据集包含了评论文本和与其相对应的七个类别的毒性得分(0-1 的实数值)。数据集分为两部分，一部分为训练集(180487 条评论)，另一部分为测试集(97320 条评论)。数据集的类别包括：toxic、severe_toxic、obsence、threat、insult、identity_hate 和正常(非毒性)。