

# LDPSA: A Large Dataset of Persian Sentiment Analysis

Ali Nazarizadeh\*

Department of Computer Engineering, Central Tehran  
Branch, Islamic Azad University  
Tehran, Iran  
computer.engineer.as@gmail.com

Taha Samavati

Iran University of Science and Technology, Faculty of  
Computer Engineering  
Tehran, Iran  
taha\_samavati@alumni.iust.ac.ir

Sina Mansouri

K. N. Toosi University of Technology, Faculty of Computer  
Engineering  
Tehran, Iran  
sina.mansouri79@email.kntu.ac.ir

Minoo Sayyadpour

Dept. Math and Computer Science  
Kharazmi University  
Tehran, Iran  
std\_minoosayyadpour@khu.ac.ir

**Abstract**—Today, data is more valuable to us than gold. When observing the environment, a substantial amount of data, particularly textual information, can be identified, tagged, prepared, and published in the form of a corpus or datasets. The primary objective of our paper is to gather, prepare, tag, and develop a vast dataset of Fidibo users' opinions regarding educational content and e-books. This dataset enables in-depth analysis of emotions and opinion mining, particularly within the educational content realm. A common flaw in nearly all similar datasets in the Farsi language is their restriction to user opinions on services and products available on online platforms. The dataset we refer to as LDPSA (A Large Dataset of Persian Sentiment Analysis) offers several advantages over comparable datasets in the Persian language. Notably, this dataset consists of 253,368 comments, each categorized into 5 classes. LDPSA represents the sole extensive Iranian dataset suitable for scrutinizing educational content and e-books. Moreover, significant insights were gleaned from data analysis. For example, during the COVID-19 pandemic, Iranian individuals dedicated more time to studying and engaging with educational platforms significantly. Nearly 80% of users expressed favorable opinions concerning the informational materials available on the Fidibo website. Users' inclination towards utilizing audio books has escalated, along with other analysis referenced in the paper.

**Keywords**—Sentiment Analysis, Persian Sentiment Analysis, Corpus, Opinion Mining, Text Mining.

## LDPSA: مجموعه داده بزرگ آنالیز احساسات فارسی

علی نظری زاده<sup>۱</sup>، طاها سماواتی<sup>۲</sup>، سینا منصوری<sup>۳</sup>، مینو صیادپور<sup>۴</sup>

<sup>۱</sup> دپارتمان مهندسی کامپیوتر دانشگاه تهران مرکز

computer.engineer.as@gmail.com

<sup>۲</sup> گروه کامپیوتر دانشگاه ایران

taha\_samavati@alumni.iust.ac.ir

<sup>۳</sup> گروه کامپیوتر دانشگاه خواجه نصیرالدین طوسی

sina.mansouri79@email.kntu.ac.ir

<sup>۴</sup> گروه ریاضی و کامپیوتر دانشگاه خوارزمی

std\_minoosayyadpour@khu.ac.ir

### چکیده

امروزه داده‌ها برای ما از طلا ارزشمندتر هستند. هنگام مشاهده محیط، مقدار قابل توجهی از داده‌ها، به ویژه اطلاعات متنی، را می‌توان شناسایی، برچسب گذاری، تهیه و در قالب مجموعه داده منتشر کرد. هدف اصلی مقاله ما جمع‌آوری، تهیه، برچسب گذاری و توسعه مجموعه وسیعی از نظرات کاربران فیدبک در مورد محتوای آموزشی و کتاب‌های الکترونیکی است. این مجموعه داده تجزیه و تحلیل عمیق احساسات و عقیده‌ها، به ویژه در حوزه محتوای آموزشی را امکان‌پذیر می‌کند. یک نقص رایج تقریباً در تمام مجموعه داده‌های مشابه در زبان فارسی، محدودیت آنها به نظرات کاربران در مورد خدمات و محصولات موجود در پلتفرم‌های آنلاین است. مجموعه داده‌ای که ما از آن به عنوان LDPSA (مجموعه داده‌های بزرگ تحلیل احساسات فارسی) یاد می‌کنیم، مزایای متعددی نسبت به مجموعه داده‌های قابل مقایسه در زبان فارسی دارد. قابل ذکر است، این مجموعه داده شامل ۲۵۳۳۶۸ نظر است که هر کدام در ۵ کلاس طبقه‌بندی می‌شوند. LDPSA تنها مجموعه داده گسترده ایرانی است که برای بررسی دقیق محتوای آموزشی و کتاب‌های الکترونیکی مناسب است. علاوه بر این، بینش‌های قابل توجهی از تجزیه و تحلیل داده‌ها به دست آمد. به عنوان مثال، در طول همه‌گیری COVID-19، افراد ایرانی زمان بیشتری را به مطالعه و تعامل با بسترهای آموزشی اختصاص دادند. نزدیک به ۸۰ درصد از کاربران نظر مساعد خود را در مورد مطالب اطلاع‌رسانی موجود در وب‌سایت فیدبک اعلام کردند. تمایل کاربران به استفاده از کتاب‌های صوتی، همراه با تحلیل‌های دیگری که در مقاله به آنها اشاره شده است، افزایش یافته است.

**کلمات کلیدی:** آنالیز احساسات، آنالیز احساسات فارسی، پیکره متنی، نظر کاوی، پردازش متن.