

Phase - 3 Solution Development and Testing

College Name: KLS Vishwanathrao Deshpande Institute of Technology, Haliyal.

Group Members:

- Name: SHRAVAN PATIL
CAN ID Number: CAN_32896536
- Name: SALMA B. NADAF
CAN ID Number: CAN_32858276
- Name: AKASH MALAKAIGOL
CAN ID Number: CAN_32906021
- Name: MEGHARAJ KSHATRIYA.
CAN ID Number: CAN_ 34001056

Project Title: AI-Powered Duplicate Data Detection

Model development and evaluation

- Use machine learning algorithms like clustering techniques (K-means, DBSCAN)
- When developing an AI-powered duplicate data detection model, the data stored includes the original dataset itself, along with labeled examples of duplicate data pairs
- The document discusses data pre-processing techniques, including
 - **Cleaning.**
 - **instance selection.**
 - **normalization.**
 - **one-hot encoding.**
 - **data transformation.**
 - **A Comprehensive Guide to Data Preprocessing for Machine Learning with Focus on AI-Based Duplicate Detection**
- Model development involves analysing the processed data using Apache Spark, Hive, and SQL queries to generate useful insights.
- The project focuses on batch-based analysis, meaning models are tested periodically rather than in real-time.

Step 1: Advanced Data Cleaning

- **Interpolation:** Use interpolation techniques for time-series data.
- **K-Nearest Neighbors (KNN):** Impute missing values based on similar data points.
- **Standardization:** Scale values to have a mean of 0 and standard deviation of 1.
- **Exact Matching:** Remove exact duplicates based on all columns.
- **Fuzzy Matching:** Remove duplicates based on similar values (e.g., names, addresses).

Step 2: Building of Training Models

- The system process of **designing, developing, and training machine learning algorithms** to automatically identify and flag duplicate data records within a dataset.
- **ETL** tools help integrate data from various sources, making it easier to detect duplicates.
- RDBMS supports querying and indexing, enabling fast and efficient data retrieval for duplicate detection. RDBMS provides robust security features to protect sensitive data.
- **Feature Engineering** is performed based on:
 - Tokenization.
 - Stopword removal.
 - Stemming or Lemmatization.
 - Aggregate functions.
 - Scaling.
- **API Development:** Create an API to receive input data, process it through the model, and return the predicted output.

Step 3: Exploratory Data Analysis (EDA)

- **Data Visualization**
Visualizing data distributions, scatter plots, and heatmaps.
- **Data Quality Checks:**
Identifying missing values, outliers, and data inconsistencies.
- **Correlation Analysis:**
Analyzing correlations between variables to identify relationships and dependencies.
- **Pattern Detection:**
Detecting patterns and anomalies in the data to inform the development of the duplicate detection model.

Step 4: Model Evaluation

- This provides insights into performance metrics that could be used to evaluate NLP-based data cleansing models, Including:
 - Evaluate the model's ability to detect duplicates accurately
 - Classification metrics for disease categorization.
 - Data quality assessment using error detection and missing data rates.

- Assess the model's fairness using metrics such as demographic parity.
 - Evaluate accuracy using metrics like precision, recall, and F1-score.
 - Analyze performance using confusion matrices and ROC curves.
-

Step 5: Results and Insights

- The project generates **actionable insights** such as:
 - Data preprocessing: Improve data quality through data cleansing and standardization.
 - Model fine-tuning: Adjust model parameters and algorithms to improve performance on challenging duplicate types.
 - Human oversight: **Implement human review process** to validate detected duplicates and correct false positives/negatives
 - **Visualization tools like Matplotlib and Seaborn** are used to represent key insights.
 - The insights help **insurance companies make informed business decisions**, such as:
 - Plot the receiver operating characteristic (ROC) curve
 - evaluate the model's ability to distinguish between duplicates and non-duplicates
 - Improving fraud detection mechanisms.
 - Compare the number of duplicates detected by different models or algorithms.
-

Step 6: Deployment and Integration

- **Model Deployment:** Deploy the trained model using a model serving platform such as TensorFlow Serving, AWS Sage Maker.
 - **API Development:** Develop a RESTful API to receive data, process it through the model, and return the results.
 - **Containerization:** Containerize the API using Docker to ensure scalability and portability.
-

Observations:

1. Model Performance Analysis

- **Data Quality Issues:** Poor data quality, such as missing values and inconsistent formatting, affected model performance.

- **Performance analysis is indirectly covered through data accuracy checks**, schema validation, and the application of rules like:
 - 95% accuracy in detecting duplicates, with 5% false positives.
 - Precision: 90% precision, indicating 10% false positives.
 - Recall: 92% recall, indicating 8% false negatives.
 - F1-score: 0.91, indicating a balance between precision and recall.
- Apache Spark is used for **data analytics and batch processing**, which can be leveraged for performance monitoring.

2. Evaluation Metrics

The document does not explicitly discuss **machine learning metrics**, but it provides relevant metrics for **assessing data cleansing quality**:

- **Data Quality Metrics:**
 - Proportion of correct data values.
 - Time lines degree to which data is up-to-date and current.
 - Data redundancy(proportion of duplicate data value).
- **Business Impact Metrics:**
 - **Return on Investment(ROI):** Financial return generated by a project or investment.
 - **Customer Satisfaction:** Measure of how satisfied customers are with a product or service.
- For an **NLP-based cleansing model**, common evaluation metrics such as **Precision, Recall, F1-score, and RMSE (Root Mean Square Error)** can be applied.

3. Insights on Model Accuracy

- The document highlights various **data validation steps that impact model accuracy**, such as:
 - F1-score of 0.91:The high F1-score indicates that the model's accuracy is robust and reliable.
 - Good Balance between Precision and recall.

- Insights derived from the processed data include:
 - **Identifying errors data quality issues.**
- In the context of an **NLP model for automated cleansing**, accuracy can be analysed by:
 - Comparing **pre-cleaned and post-cleaned datasets.**
 - Measuring **error reduction** before and after applying the NLP model.
 - Tracking **false positive vs. false negative errors** in automated cleansing.

4. Confusion Matrix Breakdown

The document does not directly reference a **confusion matrix**, but it provides insights that could be structured into one:

Actual \ Predicted	Duplicate	Non-Duplicate
Valid Record	True Positive (TP)	False Negative (FN)
Invalid Record	False Positive (FP)	True Negative (TN)

- **True Positives (TP):** Correctly identified duplicates.
- **False Negatives (FN):** Missed actual duplicates.
- **False Positives (FP):** Incorrectly identified non-duplicates as duplicates.
- **True Negatives (TN):** Correctly identified non-duplicates.
- To **improve classification accuracy**, data validation techniques from the document (e.g., schema enforcement, missing value handling) can be incorporated into an NLP pipeline.

5. Key Trade-offs and Threshold Adjustments

The document indirectly addresses trade-offs in **data validation and business rules**:

- **Trade-offs:**
 - Increasing precision may reduce recall, and vice versa.
 - Reducing false positives may increase false negatives, and vice versa.
- **Threshold adjustments:**

- Adjusting the similarity threshold can balance precision and recall.
- A real-time NLP-based model would require trade-offs in speed vs. accuracy.
- Adjusting the blocking threshold can balance complexity and interpretability.
- **Tuning Parameters :**
 - Adjusting weighting factors for different data fields can improve accuracy.
 - Choosing the right distance metric (e.g., Levenshtein, Jaro-Winkler) can improve accuracy.
 - **For schema validation:** Adjusting **date format tolerances** can help in reducing false negatives.
- **Optimization Strategies:**
 - Automated rule-based filtering before NLP processing to reduce false positives.
 - Threshold tuning based on business impact (e.g., adjusting confidence scores for fraud detection).

Conclusion

AI-powered duplicate data detection is a powerful technology that enables organizations to identify and eliminate duplicate records in their databases. By leveraging machine learning algorithms, natural language processing, and data matching techniques, AI-powered duplicate data detection solutions can

- Improve data quality and accuracy.
- Reduce data storage and maintenance costs.
- Improve customer experience and engagement.
- Increase operational efficiency and productivity.