

Discussion

MoneyLion Machine Learning Scientist Take Home Assessment

Predict Bank Transaction Category Challenge by Johnston Yap
johnstonyap.jy@gmail.com | +6011-1113 3792

- Time Spent : 13 hours

1. Data Analysis and Model Development

The datasets used in this task include `bank_transactions` and `user_profiles`. The primary dataset, `bank_transactions`, contains various attributes such as client ID, bank ID, account ID, transaction ID, transaction date, description, amount, and category. The target variable for this task is the category of the transaction, which we aim to predict based on the description.

The initial data analysis involves checking for NaN values (only `categories` have NaNs with 257 rows), understanding the distribution of transaction categories, and identifying potential issues such as typos or irrelevant categories (e.g., 'Transfer', 'Uncategorized'). Entries like 'Transfer' cannot be defined as Deposit, Credit, or Debit, and 'Uncategorized' needs to be excluded due to its lack of guidance on training and validating datasets. The data was cleaned by removing such entries and merging similar categories (e.g., 'Bank Fee' and 'Bank Fees'), addressing typos in the labels.

Reasoning and Justification for Model Architecture

Given the text-based nature of the 'description' column, which is used to predict the transaction category, the following preprocessing steps and model choices were made:

a. Text Preprocessing

- Removing extra spaces, punctuation, and numeric values.
- Tokenizing the text.
- Lemmatization and removing stop words were considered but commented out to maintain features, which increased the F1-Score.

b. Feature Extraction on Techniques for Different Models:

Polynomial Features & Interaction Terms :

- Logistic Regression (LR): Helpful for capturing non-linear relationships between features, making it useful when you suspect non-linearity.
- Support Vector Machines (SVM): Can be beneficial, especially with polynomial kernels, to handle complex relationships and create higher-dimensional spaces.
- Naive Bayes (NB): Generally not useful because NB assumes feature independence. Non-linear transformations won't affect its performance.
- Random Forests (RF): Less relevant, as RF's tree-based structure naturally captures non-linear interactions.

Feature Scaling :

- Logistic Regression (LR): Crucial, as LR's coefficients are sensitive to feature scales. Different scales can lead to poor model performance.

- Support Vector Machines (SVM): Essential for maximizing the margin between classes. Features on different scales can impact the SVM's performance.
- Naive Bayes (NB): Not critical, since NB is based on feature probabilities and doesn't depend on feature magnitudes.
- Random Forests (RF): Not necessary, as RF's tree-based nature is not affected by feature scales.
-

TF-IDF (Implemented) :

- Logistic Regression (LR): Works well, as TF-IDF converts text data into numerical features that LR can use for modeling.
- Support Vector Machines (SVM): Effective, as TF-IDF transforms text data into a format SVM can handle to find the optimal separating hyperplane.
- Naive Bayes (NB): Often used, especially with variants like Multinomial Naive Bayes, to weigh word importance in text data.
- Random Forests (RF): Can be used, though RF may be less sensitive to term frequencies. However, it can still benefit from TF-IDF features.

| Model Name | Polynomial Features & Interaction Terms | Feature Scaling | TF-IDF |
|-------------------------|--|-----------------|---|
| Logistic Regression | Useful for capturing non-linearity | Important | Suitable for numerical feature modeling |
| Support Vector Machines | Beneficial for complex relationships | Crucial | Effective for SVM classification |
| Naive Bayes | Not relevant | Not critical | Often used, especially with variants |
| Random Forests | Less relevant, captures interactions naturally | Not necessary | Can be used, but less sensitive to term frequencies |

c. Model Selection

- **Naive Bayes:** Suitable for text classification tasks due to its simplicity and efficiency.
- **Logistic Regression:** Effective for binary and multiclass classification tasks.
- **Random Forest:** A robust ensemble method that handles overfitting well and provides feature importance.
- **Support Vector Machine (SVM):** A powerful classifier for high-dimensional spaces, though computationally intensive.

d. Hyperparameter Tuning with K-Fold Cross-Validation

To rigorously evaluate the performance of the models, 5-fold cross-validation was implemented. This method divides the dataset into five distinct folds, using each fold once as a validation set while the remaining four folds are used as training data. This process is repeated five times, ensuring every data point is used for both training and validation.

Key Benefits of 5-Fold Cross-Validation:

- **Reduced Overfitting:** Assesses the model's ability to generalize to unseen data by validating on multiple subsets.
- **Performance Metrics:** Provides a more reliable estimate of model performance by averaging metrics (accuracy, F1-score, etc.) across all folds.
- **Efficient Use of Data:** Ensures all data points are used for both training and validation, leading to a more comprehensive evaluation.

2. Functional Code and Reasoning

a. Loading Data:

- Read `bank_transaction.csv` data from CSV files.
- Display initial information for understanding the structure.

b. Data Visualization:

- Plot the distribution of transaction categories.
- Display description field from each category to examine differences.

c. Data Preprocessing:

- Initial data analysis involved checking for NaN values, understanding the distribution of transaction categories, and identifying potential issues such as typos or irrelevant categories (e.g., 'Transfer', 'Uncategorized'). The data was cleaned by removing such entries and merging similar categories (e.g., 'Bank Fee' and 'Bank Fees').
- Preprocessing descriptions to clean and tokenize text.

d. Feature Extraction:

Using `TF-IDF` to transform text data into numerical features.

- TF-IDF transforms text into weighted numerical features that enhance model performance by emphasizing term relevance, which benefits LR, NB, SVM, and RF in capturing important patterns in text data.

| Model Name | Polynomial Features & Interaction Terms | Feature Scaling | TF-IDF |
|-------------------------|--|-----------------|---|
| Logistic Regression | Useful for capturing non-linearity | Important | Suitable for numerical feature modeling |
| Support Vector Machines | Beneficial for complex relationships | Crucial | Effective for SVM classification |
| Naive Bayes | Not relevant | Not critical | Often used, especially with variants |
| Random Forests | Less relevant, captures interactions naturally | Not necessary | Can be used, but less sensitive to term frequencies |

e. Data Splitting:

- Splitting the dataset into training and testing sets.

f. Model Training and Evaluation:

- Training and evaluating models (Naive Bayes, Logistic Regression, Random Forest, SVM).
- Implementing 5-Fold Cross-Validation ensures comprehensive evaluation.
- Classification Report and Confusion Matrix were used to determine Accuracy and F1-Score.

3. Model Performance and Effectiveness

Model Evaluation Results (With 5-Fold Cross Validation):

| Model | Training Time (seconds) | Accuracy (%) | F1 Score |
|---------------------|-------------------------|--------------|----------|
| Naive Bayes | 0.376576 | 79.89 | 0.78 |
| Logistic Regression | 35.14 | 88.62 | 0.88 |
| Random Forest | 201.46 | 89.93 | 0.90 |
| SVM | 365.07 | 89.51 | 0.89 |

Conclusion:

- The **Random Forest** model achieved the highest F1 score of 0.90, indicating it is the best-performing model among those evaluated.
- The **Logistic Regression** model achieved a high F1 score of 0.88 and a fast training time of 35.14 seconds, indicating it is the most efficient model among those evaluated.

6.4 Future Development Plans

Next 1 Month:

- **Technique:** Experiment with feature scaling over our TF-IDF features to enhance performance.
- **Hyperparameter Tuning:** Use grid search or random search to fine-tune model hyperparameters, alongside existing K-Fold Cross-Validation.
- **Feature Engineering:** Explore additional features like transaction amount, user profile information, and transaction date to enrich the model's input data.

Next 3 Months:

- **Feature Improvement:** Leverage data from `user_profile.csv` to enhance classification and potentially increase conversion rates.
- **Model Ensemble:** Combine predictions from different models to boost accuracy, possibly through a voting system.
- **Deep Learning Models:** Experiment with advanced models like LSTM or BERT to capture subtle information and improve F1 scores.
- **Data Feeding:** Implement a system for continuous updates with new data.
- **Deployment:** Develop a scalable API for real-time predictions and integrate it into a production environment.

Conclusion

This project focused on developing a machine learning model to predict bank transaction categories based on descriptions. Using **TF-IDF** proved to be an efficient choice for feature extraction across models. The **Logistic Regression** model demonstrated a good balance between accuracy and training time, while the **Random Forest** model performed best in terms of F1 score. Future work will involve improving model performance and preparing for practical deployment.