

# Project Proposal

Spring 2024

## Title: Credit Card Fraud Detection using Machine Learning

### 1. Objective

Our aim is to enhance credit card fraud detection methods by delving deeper into the features that define fraudulent transactions. Leveraging Feature Space Enrichment (FSE) techniques, we seek to expand the feature space of credit card transaction datasets. This project endeavors to compare a baseline fraud detection model against a Holistic Credit Card Fraud Detection (HCCFD) model derived post FSE. We'll explore two variants of the HCCFD model: HCCFD-1, trained on a single enriched credit card fraud dataset, and HCCFD-2, a pipeline consisting of two sub models - one focusing on demographic features and the other on transaction features, both corroborating each other.

### 2. Motivation

The prevalence of credit cards in modern society has led to a significant increase in credit card fraud cases, despite the integration of chip cards and existing protection systems. This issue is crucial due to the rise in online transactions and e-commerce platforms. Credit card fraud can result in substantial financial losses and damage to reputation for financial institutions and cardholder. To address this challenge, we propose a project that focuses on building a holistic credit card fraud detection model using machine learning techniques. Our model will leverage both demographic and transaction-related data to enhance the robustness of the fraud detection system against changing fraud tactics.

### 3. Related Work

Credit card fraud detection has attracted significant research attention due to its pivotal role in maintaining financial security and stability. Numerous studies have explored the utilization of machine learning techniques to address this pressing challenge like [1][2][3]. Gao et al. (2021), have delved into machine learning methods to tackle this challenge, emphasizing the need for

robust models capable of adapting to evolving fraud tactics and mitigating the impact of concept drift [1]. Hande et al. (2019) further contributed to this domain by proposing tailored machine learning approaches specifically for credit card fraud detection, showcasing the effectiveness of algorithms like decision trees and support vector machines in accurately distinguishing fraudulent transactions [2].

Israel et al. (2023) extended this research focus by emphasizing feature selection's intricate role in enhancing fraud detection system efficacy, particularly in the context of financial fraud detection, with a specific emphasis on credit card fraud. Their exploration of various feature combinations and algorithmic approaches provides valuable insights for developing sophisticated and adaptive fraud detection mechanisms [3]. Additionally, the concept drift problem discussed in [4] explains how changing data due to the dynamic nature of credit card fraud is a one of the important challenges in credit card fraud detection using machine learning apart from the challenge of data skewed towards non-fraudulent transactions i.e data imbalance [1].

## **4. Methodology**

Our methodology operates on the assumption that all datasets used in this project, sourced from various origins, represent the same set of customers. This assumption simplifies data alignment processes, facilitating the merging of the base dataset with external datasets to enable FSE. While this approach may yield machine learning models with moderate performance due to data alignment constraints, it serves as a foundational guide for holistic credit card fraud detection.

The key concept driving our methodology is the notion that every credit card transaction encompasses both demographic and transactional data. We propose to enrich demographic data using publicly available sources such as spending habits data, salary distributions, and employment types. Similarly, we'll enhance transactional data using related payment channel transaction data like app transfers and net banking transactions.

Following enrichment, subsets of data will either contribute to creating new features for the master model's final training data or participate in the creation of sub models. These sub models will collaborate to make decisions regarding transaction fraudulence, with each sub model providing its perspective on whether a transaction is fraudulent or not.

The steps involved can be elaborated as follows:

### **1. Data Preprocessing and Enrichment**

- Data Collection: Gather the credit card transaction data, including information about the transactions and the cardholders.
- Data Cleaning: Remove any errors or inconsistencies in the data, such as missing values or incorrect entries.
- Feature Enrichment: Enhance the existing data by adding more relevant information, for example, including demographic details of the cardholders and additional transaction-related data from external sources.

## **2. Baseline Model Building (HCCFD variant 1)**

- Algorithm Selection: Choose the machine learning algorithms such as decision tree, naive Bayes, logistic regression for building the fraud detection models.
- Training the Models: Teach the models to recognize patterns in the data by using the enriched credit card transaction dataset.

## **3. Building Pipeline of sub models (HCCFD variant 2)**

- Sub models: Create the individual models just like the procedure mentioned in step 2, to create demographic fraud detection model and the transaction data model.
- Corroboration: Ensure that the models work together to validate each other's findings, making the fraud detection process more robust.

## **4. Model Evaluation and Validation**

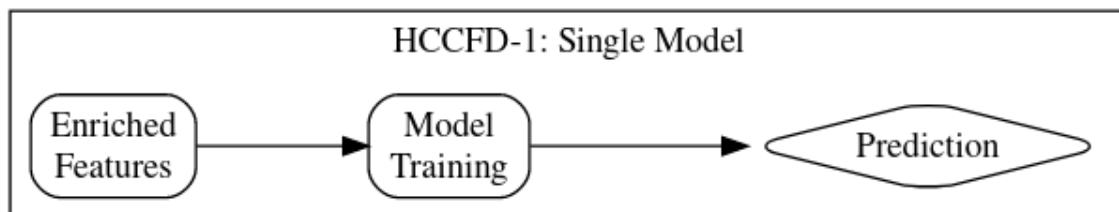
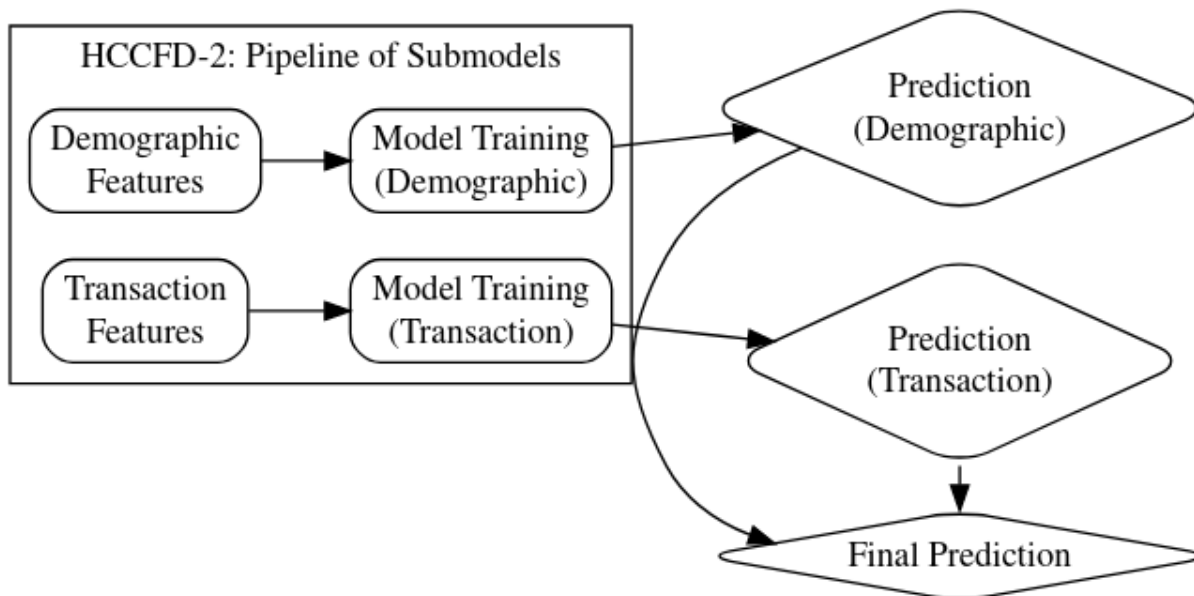
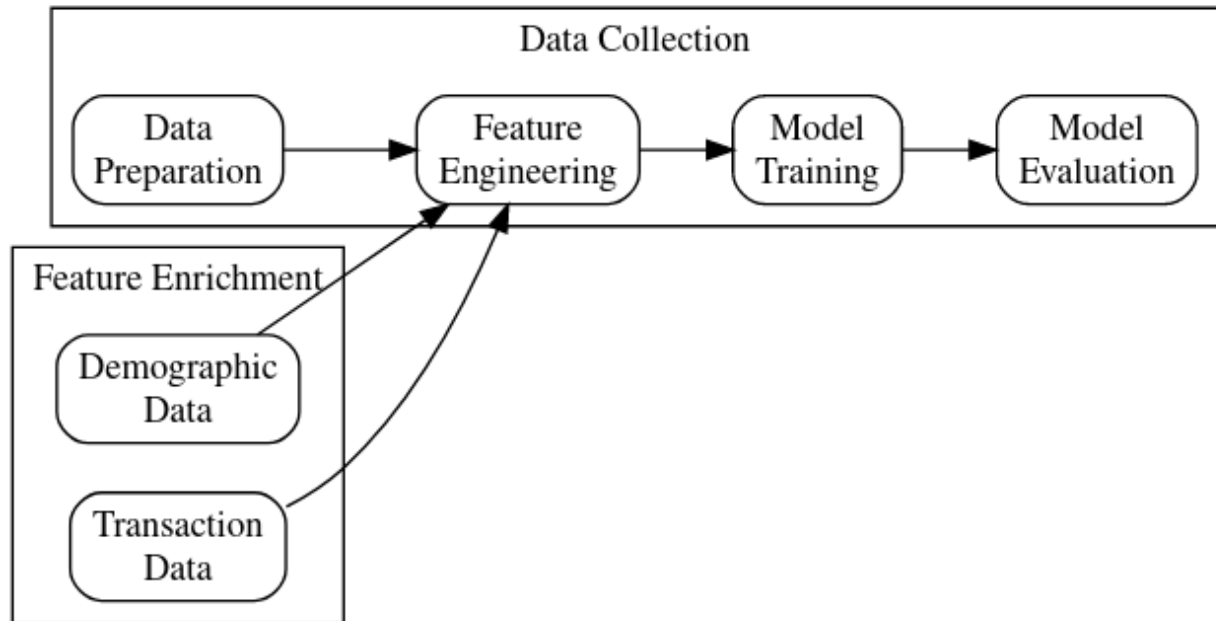
- Model Testing: Assess how well the models perform by using a portion of the data that the models haven't seen before.
- Performance Measurement: Measure the models' accuracy in identifying fraudulent transactions and their ability to avoid false alarms.

## **5. Probabilistic Analysis**

After evaluating our models' predictions, we'll conduct a probabilistic analysis to estimate the likelihood of observing specific fraudulent transaction patterns. This analysis will provide deeper insights into the confidence levels of our model predictions, aiding in decision-making regarding fraud occurrence.

## **6. Optimization Strategies**

To enhance our models' performance, we'll employ optimization strategies such as fine-tuning hyperparameters and assessing trade-offs between complexity and accuracy. By systematically optimizing our models, we aim to develop a reliable fraud detection system capable of adapting to evolving tactics while minimizing false alarms.



## 5. Deliverables:

- 1) **Codebase:** We will deliver a comprehensive and well-documented codebase encompassing all aspects of the project implementation. This includes scripts for data preprocessing, feature engineering, model training, evaluation, and deployment. The code will be organized and annotated for clarity and reproducibility, ensuring that stakeholders can easily understand and extend the implemented solutions.
- 2) **Trained Models:** The project will yield two variants of the Holistic Credit Card Fraud Detection (HCCFD) model. Each variant offers unique insights into credit card fraud detection, leveraging Feature Space Enrichment (FSE) techniques. These trained models, along with their optimized configurations, will be provided as artifacts for direct deployment or further analysis.
- 3) **General Framework:** In addition to the specific models, the project will produce a general framework outlining the conversion of specific models into holistic ones. This framework serves as a blueprint for future research and applications in fraud detection, offering guidance on integrating demographic and transactional data to enhance fraud detection systems.
- 4) **Technical Report:** A detailed technical report will be produced, outlining the project's methodology, data analysis, feature enrichment techniques, model selection, training procedures, evaluation metrics, and findings. The report will include visualizations, statistical analyses, and a discussion of the implications of the results, offering insights for both academic and practical audiences.
- 5) **Presentation Slides:** We will create a set of presentation slides suitable for academic and professional settings. These slides will cover the project's background, objectives, methodology, results, and potential applications, providing a concise overview of the project for diverse audiences.
- 6) **Documentation:** Thorough documentation accompanying the codebase will be provided, offering insights into the project's structure, dependencies, and instructions for running and deploying the models. This documentation will serve as a guide for replicating the study and implementing the predictive models in real-world scenarios, ensuring ease of use and reproducibility.

Overall, these deliverables aim to provide a comprehensive package that enables both the replication of the study, and the practical implementation of the fraud detection models in real-world settings. Additionally, the general framework offers a roadmap for future advancements in the field of fraud detection.

## 6. Resources

The successful execution of this project relies on access to the following resources:

- **Hardware:** A personal computer or access to a computing environment with sufficient computational resources to train and evaluate machine learning models. A machine with a multicore CPU.
- **Data Access:** Access to credit card transaction datasets with fraud labels is crucial for model development and evaluation. We will require access to datasets containing labeled credit card transactions, such as the Kaggle Credit Card Fraud Detection dataset.
- **Software Resources:** For Software, Jupyter Notebooks, Python and libraries such as NumPy, pandas, and scikit-learn will be employed for data analysis and model development. Furthermore, we might require Matplotlib and Seaborn for data visualization along with LaTeX for document preparation.
- **Time Allocation:** Adequate time allocation for project execution and experimentation is necessary to ensure thorough exploration of different methodologies, model architectures, and optimization strategies. This includes time for data preprocessing, feature engineering, model training, evaluation, and documentation.
- **Internet Access:** Reliable internet access is required for literature review, accessing online resources, and potential collaboration with external datasets or research findings.

All the mentioned resources are available and accessible, ensuring the feasibility of the project within a typical academic setting.

## 7. Impact

The proposed project on enhancing credit card fraud detection methods has the potential to make significant contributions to both academic research and real-world applications:

- **Improved Effectiveness:** By leveraging Feature Space Enrichment (FSE) techniques and a holistic approach, the project aims to significantly enhance the effectiveness of credit card fraud detection methods. This could lead to a reduction in financial losses and reputational damage for financial institutions and cardholders.

- **Robust Fraud Detection Systems:** Incorporating demographic and transaction-related data into fraud detection systems can enhance their robustness against evolving fraud tactics. The holistic approach proposed in this project ensures a comprehensive understanding of fraudulent patterns, thus enabling better detection and mitigation strategies.
- **Insights into Feature Space:** The project seeks to provide deeper insights into the feature space of credit card fraud datasets. Understanding the intricate relationship between various features and fraudulent activities can lead to advancements in fraud detection research, potentially uncovering new patterns and trends.
- **Practical Applications:** The developed models and methodologies are not only relevant for academic research but also hold significant practical value. Financial institutions and other stakeholders can directly benefit from the implementation of more effective fraud detection systems, thereby safeguarding their assets and maintaining trust with customers.
- **Adaptability and Scalability:** The project's framework and models can be adapted and scaled to different datasets and fraud detection scenarios beyond credit card fraud. This adaptability ensures that the research findings have broader applicability across various domains, further enhancing their impact.
- **Contributions to Machine Learning Research:** The exploration of machine learning algorithms and optimization strategies within the context of fraud detection contributes to the broader scientific understanding of these techniques. Comparative analyses and performance evaluations provide valuable insights into the strengths and limitations of different approaches, advancing the field of machine learning in fraud detection.

Overall, the proposed project not only addresses the pressing challenge of credit card fraud but also lays the groundwork for future research and practical applications in fraud detection and prevention. Its impact extends beyond theoretical advancements to tangible improvements in financial security and risk management practices.

## 8. Milestones

The successful completion of the project will involve achieving several key milestones. The proposed timeline for these milestones is outlined below:

1. **Literature Review and Related Work Analysis:** Conduct an extensive review of literature and related works on credit card fraud detection methods, machine learning

techniques, and feature space enrichment. Gain insights from previous research to inform project methodology and approach.

Target Date: 01/2024

- 2. Data Collection and Preprocessing:** Gather credit card transaction datasets and relevant demographic data. Perform data cleaning to remove errors or inconsistencies. Enrich the dataset by incorporating additional demographic and transaction-related features.

Target Date: 02/2024

- 3. Model Development - Baseline Model (HCCFD variant 1):** Select appropriate machine learning algorithms such as decision trees, naive Bayes, and logistic regression. Train the baseline fraud detection model using the enriched credit card transaction dataset.

Target Date: 02/2024

- 4. Model Development - Pipeline of Sub Models (HCCFD variant 2):** Develop individual sub models focusing on demographic features and transaction features separately. Ensure collaboration between sub models to validate each other's findings and enhance overall fraud detection accuracy.

Target Date: 03/2024

- 5. Model Evaluation and Validation:** Evaluate the performance of both HCCFD variants using a separate portion of the dataset that the models haven't seen before. Measure accuracy in identifying fraudulent transactions and assess the ability to minimize false alarms.

Target Date: 03/2024

- 6. Probabilistic Analysis:** Conduct probabilistic analysis to estimate the likelihood of observing specific fraudulent transaction patterns. Gain deeper insights into the confidence levels of model predictions and aid decision-making processes regarding fraud occurrence.

Target Date: 03/2024

- 7. Report Writing and Documentation:** Prepare a detailed technical report outlining the project's methodology, data analysis, feature enrichment techniques, model development, evaluation metrics, and findings. Document the codebase comprehensively for replication and future reference.

Target Date: 03/2024



- 8. Presentation Preparation:** Create presentation slides summarizing the project background, objectives, methodology, results, and potential applications. Tailor presentations for academic and professional settings to effectively communicate project insights.

Target Date: 04/2024

- 9. Final Deliverables Submission:** Submit the final project report, presentation materials, codebase, trained models, and general framework documentation. Ensure all deliverables are organized and accessible for stakeholders.

Target Date: 04/2024

The milestones outlined above are subject to adjustment based on project progress, resource availability, and unforeseen challenges encountered during the implementation phase.

## 9. References

1. Y. Gao, S. Zhang, and J. Lu, "Machine Learning for Credit Card Fraud Detection," in Proceedings of the 2021 1st International Conference on Control and Intelligent Robotics (ICCIR '21), New York, NY, USA: Association for Computing Machinery, 2021, pp. 213-219. DOI: [10.1145/3473714.3473749](https://doi.org/10.1145/3473714.3473749).
2. T. Hande et al., "Credit Card Fraud Detection using Machine Learning," Int. J. Eng. Adv. Technol., 2019. Available: <https://www.ijert.org/research/credit-card-fraud-detection-using-machine-learning-and-data-science-IJERTV8IS090031.pdf>
3. O. Israel et al., "Financial Fraud Detection using Machine Learning: Credit Card Fraud," Int. J. Recent Eng. Sci., 2023. Available: <https://www.ijresonline.com/assets/year/volume-10-issue-3/IJRES-V10I3P104.pdf>
4. O. S. Adebayo et al., "Comparative Review of Credit Card Fraud Detection using Machine Learning and Concept Drift Techniques," Int. J. Comput. Sci. Mobile Comput., 2023. Available: <https://ijcsmc.com/docs/papers/July2023/V12I7202310.pdf>