

# Capstone Project Proposal



Jessica Ofoh

## Business Goals

<b>Project Overview and Goal</b>  What is the industry problem you are trying to solve? Why use ML/AI in solving this task? Be as specific as you can when describing how ML/AI can provide value. For example, if you're labeling images, how will this help the business?	<p>Project Name: Fraudulent Bank Transaction Detection</p> <p>The industry problem that I aim to solve under this project is the automatic detection of fraudulent bank transactions in the finance sector. With the increase in accuracy of ML and DL models in the cybersecurity space, coupled with astonishing fast computational speed, I believe that AI/ML can be assuredly implemented in this project to make it feasible and practical in the real world. Moreover, as the project helps in labeling transactions as safe or fraudulent, it would critically help financial institutions to stop malicious activities on their platforms by a huge factor.</p>
<b>Business Case</b>  Why is this an important problem to solve? Make a case for building this product in terms of its impact on recurring revenue, market share, customer happiness and/or other drivers of business success.	<p>Several million fraudulent transactions take place every year in various financial institutions that become liable to them in terms of revenue lost in rectifying the damage imposed due to the malicious activities. Moreover, the fact that there's no concrete procedure to detect such activities, has made the people who execute these fraudulent transactions become confident of their capabilities and get motivated to keep executing frauds. Therefore, the creation of a system that can detect such fraudulent transactions is the need of the hour for financial institutions.</p>
<b>Application of ML/AI</b>	<p>The precise task in the implementation of the fraudulent bank transaction detection system, would be to LABEL the incoming transaction requests by the customers onto their platform to be SAFE or FRAUDULENT. In case the</p>

What precise task will you use ML/AI to accomplish? What business outcome or objective will you achieve?

transaction is termed as fraudulent, the customer's account will be locked up (denial of services by institution, cautionary emails and messages sent to customer among other steps) until there's concrete information or proof on the safety of the customer's account. This would drastically decrease the amount of fraudulent activities on financial institutions' platform(s).

## Success Metrics

### Success Metrics

What business metrics will you apply to determine the success of your product? Good metrics are clearly defined and easily measurable. Specify how you will establish a baseline value to provide a point of comparison.

The business metric that I would utilize to determine the success of the product, would be the (decrease in) amount of revenue lost due to rectification of fraudulent transactions incurred by the financial institution. If the amount of successful fraudulent cases detected goes up or the expenditure on the rectification goes down; then it would be a success for the product in general. Moreover, to establish a baseline value to provide a point of comparison for the product, we could have a comparative look at the statistics of the fraudulent cases before and after the deployment of the product.

## Data

### Data Acquisition

The acquisition of data for the fraudulent bank transaction detection can be accomplished in various methods, that include using of publically available academic data (on platforms like UCI, Kaggle, etc.), using the pre-stored and directly available historical data of the financial institution's fraudulent cases or pay for datasets on commercial data provider sites. I would

<p>Where will you source your data from? What is the cost to acquire these data? Are there any personally identifying information (PII) or data sensitivity issues you will need to overcome? Will data become available on an ongoing basis, or will you acquire a large batch of data that will need to be refreshed?</p>	<p>recommend the usage of own historical dataset coupled with publicly available dataset for the initial stages of product life cycle (and scale it up with paid data provider services as per the necessity in case accuracy or hit ratio are not satisfactory). Moreover, we would be required to take special care about PPI issues in case of utilizing the publicly available dataset, and take appropriate permission and steps to avoid issues in the future. There might also be a need to annotate several cases into safe or fraudulent for effective data usage; that might cost \$2000 – \$2500 per 1000 annotations. Furthermore, the data would be dynamically updated on a monthly or quarterly basis in order to have constantly updated models in the deployed state.</p>
<p><b>Data Source</b></p> <p>Consider the size and source of your data; what biases are built into the data and how might the data be improved?</p>	<p>The source for the data acquisition process would be an amalgamation of publically available academic dataset and institution owned historical data (for fraudulent cases). The usage of publicly available dataset can impose a bias into the data which could be the specificity of the location where the malicious activities might occur. Therefore, there must be an equilibrium in both the sources to minimize such bias to the maximum possible extent.</p>
<p><b>Choice of Data Labels</b></p> <p>What labels did you decide to add to your data? And why did you decide on these labels versus any other option?</p>	<p>The decision for choosing the data labels for this specific project is fairly simple, this is a classification problem. Therefore, the labels decided for the fraudulent bank transaction detection system are 'safe', 'fraudulent' and 'unknown' (to account for ambiguity); as this would suffice and cater to the needs of the detection process efficiently. This is a better option as the detection of fraudulent cases is most effective when the cases are labeled into a category, and thus, the aforementioned labels are chosen by considering the product as a classification system.</p>

# Model

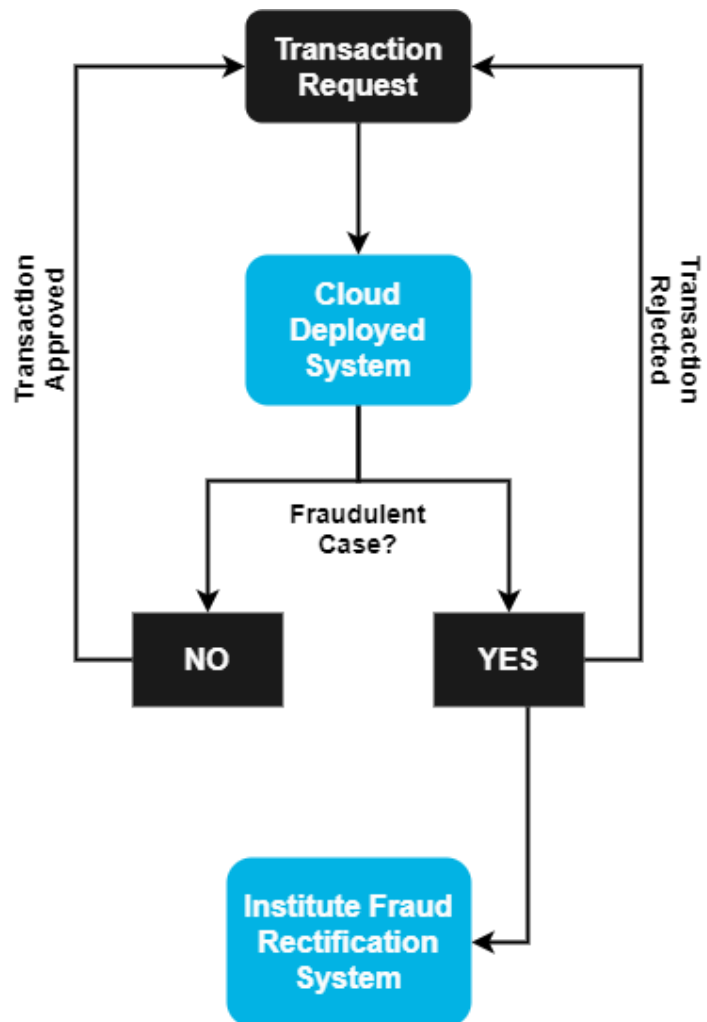
<b>Model Building</b>  How will you resource building the model that you need? Will you outsource model training and/or hosting to an external platform, or will you build the model using an in-house team, and why?	<p>I would be initially using cloud services provided by Google Cloud Platform (GCP), Amazon Web Services (AWS) or Azure to create a classification model that segregates transactions into a certain category (by labeling them) for creating a prototype. Once a certain amount of satisfaction is achieved in one of the platforms (for a certain ML model), I would be scaling up the ML model by introducing the entire dataset and creating an actual deployment level ML model. Once created, it would be deployed on financial institution's platforms. Moreover, in case the satisfaction levels decrease after a certain point in the deployment phase, I would be building a custom ML model using python and Tensorflow (or PyTorch or SKLearn) to achieve the feasible accuracy or hit ratio again.</p>
<b>Evaluating Results</b>  Which model performance metrics are appropriate to measure the success of your model? What level of performance is required?	<p>For classification models, several metrics including accuracy, hit ratio, precision, recall and F1 score can be used for effective measurement of model performance. However, in this scenario, I would be putting more emphasis on the use of hit ratio, as it clearly identifies the ratio when the model was clearly able to classify or label a transaction case appropriately. Therefore, I would be using hit ratio as the core metric (which I'd try to keep as high as possible, or above 0.98).</p>

# Minimum Viable Product (MVP)

## Design

What does your minimum viable product look like? Include sketches of your product.

As the product is mainly concerned with the classification (or segregation) of transaction requests using the ML model into safe or fraudulent cases, it is basically a backend system. Therefore, there's no UI to be portrayed here. However, the systematic workflow of the backend system can be feasibly understood via the flowchart design provided below.



<p><b>Use Cases</b></p> <p>What persona are you designing for? Can you describe the major epic-level use cases your product addresses? How will users access this product?</p>	<p>The core personas that the ‘fraudulent bank transaction detection system’ would be financial institutions from all facets of industries (governments, private banks, hedge funds, international finance organizations), as mentioned previously. Moreover, the model can also be utilized in the detection of fraudulent transactions of any kind with minor tweaks in the dataset. Hence, the product covers a significant portion of online transaction space.</p>
<p><b>Roll-out</b></p> <p>How will this be adopted? What does the go-to-market plan look like?</p>	<p>The pre-launch phase of the deployment would include the market research, initial testing, marketing and making the product deployment ready. This would suffice as the preliminary stage of the product rollout. Moreover, once pre-launch steps are taken care of, the product will be launched onto financial institutions’ platforms for active fraud detection (which would need constant and close attention in the first few months), with consistent bug fixes and relevant updates as per the performance.</p>

## Post-MVP-Deployment

<p><b>Designing for Longevity</b></p> <p>How might you improve your product in the long-term? How might real-world data be different from the training data? How will your product learn from new data? How might you employ A/B testing to improve your product?</p>	<p>The ML model designed for the ‘fraudulent bank transaction detection system’ would take the long-term perspective of model relevance to the market into consideration, by dynamically feeding the newer data (from both in-house historical dataset and publicly available dataset) on a monthly or quarterly basis as discussed earlier. This would incorporate all the updates required to be made into the system in lieu of time for staying relevant to the market and maintaining the system’s performance. Moreover, the dynamic (continuous) feeding of the data can also be cross-verified with the performance of previous and current models to deploy the model that performs better.</p>
---	--

<b>Monitor Bias</b>  How do you plan to monitor or mitigate unwanted bias in your model?	<p>The mitigation of bias, once monitored in the 'fraudulent bank transaction detection system', can be undertaken by incorporating the dataset(s) that rectify the bias in question. For instance, if the model has poor performance in detection of fraudulent cases in India, the rectification for it would be addition of dataset(s) of India origin. Therefore, the core strategy of reducing or removing unwanted bias in the system, would be the incorporation of dataset(s) that would eradicate the bias or at the very least, reduce it drastically.</p>