

# Proposal

## 1. Data set selection:

I am using a data set I found on Kaggle that is on credit card customers of a business. It consists of 10,000 customers mentioning their age, salary, marital status and more (nearly 18 features) and whether they attrited from using the card. I found this interesting because of its potential to reveal patterns associated with credit card customer churn (churn means customers quit using the card. Intuitively I would believe it's mostly based on customers' financial ability but the richness of this data set will allow me to investigate other factors potentially influencing cardholders to keep their cards or not. <https://www.kaggle.com/datasets/whenamancodes/credit-card-customers-prediction>

## 2. Objective:

What are the potential indicators or patterns for credit card customers of this firm to churn? What characteristics connect customers who churned? And what characteristics do customers who stick with the card share? By examining the network structure and attributes of influential nodes, this project aims to identify factors contributing to potential churn, allowing for the development of targeted retention strategies.

## 3. Process outline:

1. Cleaning the dataset and constructing the graph of the customer network. (nodes=customers; edges: defined based on having a large number of shared characteristics between customers)
2. Calculate **closeness centrality** and analyze the results of centrality calculations:
  - a. Identify nodes (aka customers) with high centrality
  - b. find prevalent shared characteristics among them (as potential reasons for churn or not churn)
    - i. Testing: printing out shared characteristics of a subset of high centrality nodes and verifying them manually
3. Examining the centrality of nodes in two separate groups: churn or not churn
  - a. Split the data set into two groups: churn vs not churn
    - i. Testing: verify with the known churn percentage (16.07% churn)
  - b. Identify nodes with high centrality inside each group, look at the prevalent characteristics of high centrality nodes in each group
    - i. Testing: printing out shared characteristics of a subset of high centrality nodes and verifying them manually
  - c. Compare the prevalent characteristics of both groups; use descriptive statistics (ex: mean, mode, standard deviation) to summarize the distribution of each characteristic in both groups
4. Visualize the distribution of key characteristics (and highlight key differences) for high-centrality nodes in each group

I can stop and conclude the possible characteristics that set churn and not churn customers apart, or

5. (maybe, depending on time) conduct **correlational analysis** to identify characteristics that have relatively stronger relationships with customer churn
  - a. Calculate correlation coefficients between identified key characteristics and churn status, and compare the strength of each correlation

- b. This will yield a clearer picture of how important each characteristic is in understanding potential customer churn (I need feedback on whether this is necessary or is it sufficient to stop at 4?)

## Rubric

1. Does the code run and produce reasonable/correct output (-10%)
  2. Is the dataset on which it is operating of a reasonable size (1000 vertices+) (-10%) x
  3. Is the implementation of good complexity (150+ lines of code, split in modules, not including tests) (-10%) x
  4. Does the code have tests (-10%) x
  5. Is there a good write-up describing what the project does, how to run it, what the output looks like etc. (-10%) x
  6. What is the quality of the coding (variable name selection, split of functionality in reusable functions, good use of iterators and language features (i.e. enums, structs, methods, ect) (-10%) x
- Is there a good write-up describing what the project does, how to run it, and what the output looks like? (-10%)

# Write-up

## Overview

The project analyzes the prominent characteristic categories associated with churned customers and customers who decide to continue using a bank's credit card. It operates on a data set consisting of 10,000+ customers, including their age, salary, marital status, and more (nearly 12 features), and indicates whether they have stopped using the card. To run it, simply run `cargo run --release` on Terminal.

## Data Analysis Process:

1. **Graph Construction:** The project begins by reading the data and constructing an undirected graph, where nodes represent customers and connections indicate shared characteristics.
2. **Customer Classification:** Nodes are categorized into churn and not churn groups, and high centrality nodes in both groups are identified using functions from the `graph_utils` module.
3. **Shared Characteristic Identification:** The top shared characteristic categories and their distribution are printed using functions from the `customer` module.

## Modules Overview:

1. **Main.rs:**
  - a. The main function reads the CSV file, extracts information into the `Customer Struct` (defined in `Customers.rs`), and constructs an undirected graph using the `construct_graph` function from `Graph_utils.rs`. Nodes represent customers, connected when sharing characteristics. The undirected graph is chosen for its suitability in depicting shared characteristics without directional implications.
  - b. The `Customer Struct` is then divided into churned and not churned groups. Centrality is calculated for each node in both groups, identifying nodes with high centrality through the `calculate Centrality` and `identify_high Centrality nodes` functions from `Graph_utils.rs`.
  - c. Finally, the `print_top_shared_characteristics` function from `Customers.rs` displays the top 4 shared characteristic categories and their characteristic distributions for each group.
2. **Graph\_utils.rs:**
  - a. This module handles graph construction and analysis. It imports `petgraph` and defines functions such as `construct_graph` that iterates through customer nodes to determine neighbors using `determine_neighbor` based on shared characteristics. Centrality is calculated using the `dijkstra` function, and `identify_high Centrality nodes` filter nodes with centrality scores exceeding the average times a threshold factor.
  - b. The threshold factor, currently set at 1.1 in `main.rs`, can be adjusted for dataset size and concentration. Lowering it may yield more meaningful results if there's minimal output.
3. **Customers.rs:**
  - a. Defines the `Customer Struct` detailing characteristics, including an embedded struct named `OneHotEncoding` for categorical variables, imported in `main.rs`. Functions in this module, such as `print_top_shared_characteristics`, process high-centrality nodes. `find_top_shared_characteristics` iterates through neighbors using `get_shared_characteristics` to record and return the top 4 shared characteristics for each node-neighbor pair.

- b. Results are printed and formatted with counts representing shared characteristics, organized by characteristic (e.g., Married, Single, etc.), and further categorized (e.g., Marital Status, Card Type).

```

Churn High Centrality Nodes
Prevalent characteristic categories and their compositions:
Number of Products Purchased, (Total Count: 12 - 2%)
  2: 12 (100%)
Card Type, (Total Count: 144 - 24.3%)
  Blue: 144 (100%)
Education Level, (Total Count: 10 - 1.7%)
  Graduate: 10 (100%)
Age, (Total Count: 182 - 30.7%)
  2: 182 (100%)
Month inactive, (Total Count: 48 - 8.1%)
  3: 48 (100%)
Number of Contacts from Bank (past 12 months), (Total Count: 64 - 10.8%)
  3: 64 (100%)
Marital Status, (Total Count: 132 - 22.3%)
  Married: 132 (100%)

Not Churn High Centrality Nodes:
Prevalent characteristic categories and their compositions:
Education Level, (Total Count: 1980 - 1.8%)
  High School: 74 (3.7%)
  Graduate: 1906 (96.3%)
Month inactive, (Total Count: 6888 - 6.3%)
  3: 1100 (16%)
  2: 5116 (74.3%)
  1: 672 (9.8%)
Income Range, (Total Count: 70 - 0.1%)
  $60K - $80K: 70 (100%)
Number of Products Purchased, (Total Count: 162 - 0.1%)
  4: 84 (51.9%)
  6: 78 (48.1%)
Age, (Total Count: 41146 - 37.4%)
  2: 41146 (100%)
Number of Contacts from Bank (past 12 months), (Total Count: 9528 - 8.7%)
  2: 6250 (65.6%)
  3: 3278 (34.4%)
Card Type, (Total Count: 35744 - 32.5%)
  Blue: 35744 (100%)
Marital Status, (Total Count: 14486 - 13.2%)
  Married: 11830 (81.7%)
  Single: 2656 (18.3%)

```

- c. Some categories internally segregate characteristics using ranges, with each range accompanied by a label for better interpretation.

Category	Labels for Each Characteristic				
	1	2	3	4	5
	Content of each characteristic Label				
<b>Age</b>	20-30	30-40	40-50	>50	
<b>Mon_w_bank</b> (how many months the customer stayed as a customer)	20-30	30-40	40-50	>50	
<b>Transactions_amount</b> (average dollar amount of transactions on card)	500<	500-1000	1000-1500	1500-2000	>2000
<b>Num_transactions</b> (average number of transactions on card)	<10	10-20	20-30	30-40	>40
<b>Avg_card_utilize</b> (Average Card Utilization Ratio (divide your balance by your credit limit))	<0.1	0.1-0.2	0.2-0.3	0.3-0.4	>0.4

- d. Customer.rs also includes tests. Within the test module, two simulated customers are generated, and the functions `test_shared_characteristics` and `test_determine_neighbor` are employed to assess the functionality of `shared_characteristics` and `determine_neighbor`.

### **Output Analysis:**

The output aligns with our intuition. One noteworthy shared characteristic is the quantity of products customers acquire from the bank. Among customers retaining the card, purchasing more products is a common shared characteristic. This observation is logical, as customers who have already bought numerous items from the bank are intuitively more inclined to keep using the credit card from the same institution. Additionally, a greater percentage of shared characteristics among customers with high centrality in the retaining card group is related to having a higher income. This observation is consistent with the intuition that individuals with more financial resources are more likely to persist in using a credit card.

### **Citation:**

- Data source: Kaggle
  - <https://www.kaggle.com/datasets/whenamancodes/credit-card-customers-prediction>
- Petgraph and Dijkstra reference: petgraph documentation
  - <https://docs.rs/petgraph/latest/petgraph/>
- Centrality formulas: Wikipedia, Towards Data Science
  - <https://en.wikipedia.org/wiki/Centrality>
  - <https://towardsdatascience.com/notes-on-graph-theory-centrality-measurements-e37d2e49550a>
- Concept clarification, debugging, write-up grammar revision and comment help: ChatGPT
  - <https://chat.openai.com/>
- Credits: Penny Lin (data set selection), TingYi Wu (debugging support)