

DSC630-T301 Predictive Analytics (2223-1)

Professor Fadi Alsaleem

Milestone 4, Project Proposal

01/29/2022

Each phase of the process:

- 1. Business understanding
  - A. Assess the Current Situation
    - a. Inventory of resources
    - b. Requirements, assumptions and constraints
    - c. Risks and contingencies
  - d. Terminology
    - e. Costs and benefits
  - B. What are the Desired Outputs
  - C. What Questions Are We Trying to Answer?
- 2. Data Understanding
  - A. Initial Data Report
  - B. Describe Data
  - C. Initial Data Exploration
  - D. Verify Data Quality
    - a. Missing Data
    - b. Outliers
  - E. Data Quality Report
- 3. Data Preparation
  - A. Select Your Data
  - B. Clean the Data
    - a. Label Encoding
    - b. Drop Unnecessary Columns
    - c. Altering Datatypes
    - d. Dealing With Zeros
  - C. Construct Required Data
  - D. Integrate Data
- 4. Exploratory Data Analysis
- 5. Modelling
  - A. Modelling Technique
  - B. Modelling Assumptions
  - C. Build Model
  - D. Assess Model
- 6. Evaluation
- 7. Deployment

Template Source : <https://www.sv-europe.com/crisp-dm-methodology/>

Abstract:

-This term-end project aims to evaluate the Students should be able to identify a business problem to address through predictive analytics. The goal is to select appropriate models and model specifications and apply the respective methods to enhance data-driven decision-making related to the business problem.

Students will identify the potential use of the process of analytics, formulate the problem, identify the right sources of data, analyze data and prescribe actions to improve not only the predictive decision making but also the outcome of decisions.

1. Stage One - Determine Business Objectives and Assess the Situation

1.1 Assess the Current Situation

We are targeting a sub-set of survey respondents who we have identified as gainfully and regularly employed and potentially in need of financial services such as savings accounts or CDs.

Business Application : The goal of this project is to assist in the marketing of financial services to customers based on self identified demographic data from the Bureau of Labor Statistics, allowing for the creation of pools of demographically sorted candidates for client services. Based on our findings, we can craft tailored survey data to better target ideal customers for our clients. This goal is met by creating a predictive model to determine if a given group of respondents is likely to be a good fit for our clients' financial products and as a result an ideal target for marketing campaigns attempting to solicit customers for those products.

1.1.1. Inventory of resources

List the resources available to the project including:

- Personnel: Prashant, Jay
- Data: ATUS 2019's Personal Data
- Computing resources: Personal computers
- Software: Jupyter Notebook(Jupyter), R Studio, IBM SPSS

1.1.2. Requirements, assumptions and constraints -

- Major Requirements: A valid predictive variable to determine viability of a candidate based on survey response in order to validate the output of our predictive model that would then enable us to target specific demographics based on respondent inputs.
- A pool of survey respondents providing demographic data about their employment status as well as details about their willingness to work and relevant experience (previous employment related responses).
- Major Constraints: A possible constraint is a significant number of declined responses from those surveyed. Without enough input, we won't be able to make a meaningful prediction based on the data, or lack thereof.
- The model also relies on a candidate's real time behavior, rather than a prepared list of responses. Because of this, there is a need for respondents to respond to the survey candidly and truthfully.

1.1.3.Risks and contingencies

- Risks include potential PII issues with respondents, which can largely be mitigated by anonymized respondent IDs. As long as a particular individual is never personally identified the risk of PII information leaking is mitigated.
- Potential for respondents to misrepresent their individual situations in their survey replies. Survey responses are to some extent unverifiable. For example, if a respondent indicates that the reason they remain unemployed is a long term medical issue, we can't verify that claim and that could potentially skew results.
- Lack of suitable candidates for a given application based on survey responses. This can be mitigated by the model indicating viability within the candidate pool.

1.1.4.Terminology

CRISP-DM is The Cross-Industry Standard Process for Data Mining – CRISP-DM is a model of a data mining process used to solve problems by experts. The model identifies the different steps in implementing a data mining project, as described below. The model proceeds through the following steps : Business Understanding – to understand the rules and business objectives of the company. • Understanding Data – to collect and describe data. • Data Preparation – to prepare data for import into the software. • Modeling – to select the modeling technique to be used. • Evaluation – Evaluate the process to see if the technique solves modeling and creating rules. • Deployment – to deploy the system and train its users

1.1.5.Costs and benefits

- Costs:
- Primary cost associated with this project is the time of the people working on it.
- Computing resources for modeling
- Benefits:
- Social benefit of assisting the process of filling employment gaps, and potentially enabling otherwise unemployable people to find employment
- Financial benefit to the company from the ability to sell this product to the target consumer.

1.2 What are the desired outputs of the project?

Business success criteria

- Provide this information to financial institutions for targeting for financial products.

Data mining success criteria

- Finding predictive demographic data from the respondent pool in order to determine what variables are best used for candidate selection.

Produce project plan

- Adhere to CRISP-DM methodology, Data prep Understanding > Data Prep > Modeling > Evaluation > Deployment

progress.png

1.3 What Questions Are We Trying To Answer?

Can we identify survey respondent segments that are candidates for potential packaging as demographic data to sell to our customers? What traits make an individual most suitable for the financial services for our customers?

Which of the survey respondents are good fits for the financial products of our customers?

2. Stage Two - Data Understanding

The second stage of the CRISP-DM process requires you to acquire the data listed in the project resources. This initial collection includes data loading, if this is necessary for data understanding. For example, if you use a specific tool for data understanding, it makes perfect sense to load your data into this tool. If you acquire multiple data sources then you need to consider how and when you're going to integrate these.

- Our Data set is a survey dataset and it contains the question related to two sections,
- The first section is income/employment
- The second section is related to time use

2.1 Initial Data Report

Initial data collection report - List the data sources acquired together with their locations, the methods used to acquire them and any problems encountered. Report problems you encountered and any resolutions achieved. This will help both with future replication of this project and with the execution of similar future projects.

```
In [1]: # Import Libraries Required
import pandas as pd
import matplotlib.pyplot as plt
import matplotlib inline
import seaborn as sns
pd.set_option('display.max_columns', None)
from scipy.cluster.hierarchy import dendrogram
from sklearn.metrics import adjusted_rand_score
from sklearn.preprocessing import StandardScaler, normalize
from sklearn.metrics import silhouette_score
import scipy.cluster.hierarchy as shc
from sklearn.decomposition import PCA
# suppress warnings
import warnings
import seaborn as sns
from matplotlib.pyplot import subplots
from sklearn import preprocessing
from sklearn.decomposition import PCA
```

```
In [2]: #Data source:
#Source Query Location:
path = '2019 Income Job'
# reads the data from the file - denotes as CSV, it has no header, sets column headers
df = pd.read_csv(path, sep=',')
df.info # df.copy()
```

2.2 Description - Data

Data description report - ATUS Respondent File This file contains case-specific variables collected in ATUS (that is, variables for which there is one value for each respondent). These include, for example, labor force and earnings information, total time providing secondary childcare, and ATUS statistical weights.

There is one record for each ATUS Respondent.

Below is a simplified example. The variable TUCASEID identifies each household, and the variable TUFUNENO identifies each individual within the household. The example contains responses from five individuals; note that the variable TUFUNENO always has TUFUNENO = 1. In the example, each respondent has a corresponding statistical weight for use in generating estimates representative of the U.S. civilian, noninstitutionalized population (TUFUNWGT), and values for school enrollment (TESCHENR, labor force status (TELF5), and total number of minutes spent alone on the diary day (TRTALONE). The actual ATUS Respondent file contains more variables and records.

TUCASEID	TUFUNENO	TUYEAR	TUMONTH	TEABSR5N	TEERN	TEERNH10	TEERNH2	TEERNHRO	TEERNHRY	TEER
20060101020210	1	22261358.19	1	1	1	40				
20060101020211	1	5019645.31	1	1	1	350				
20060101020212	1	2926068.74	1	5	0					
20060101020213	1	2578054.07	2	5	556					
20060101020214	1	3414645.94	1	4	100					

Income & Employment:

TEERNHRY "Edited: hourly-non-hourly status" / TEERNPER "Edited: for your main job, what is the easiest way for you to report your total earnings before taxes or other deductions: hour" / TEERNNUOT "Edited: do you usually receive overtime pay, tips, or commissions at your main job?" / TEHRUSL1 "Edited: how many hours or weeks do you usually work at your main job?" / TEHRUSL2 "Edited: total hours usually worked per week (sum of TEHRUSL1 and TEHRUSL2)" / TRTALONE "Edited: individual day of worker code (main job)" / TELF5 "Edited: labor force status" / TEIOICD "Edited: industry code (main job)" / TEIOICD "Edited: industry code (main job)" / TELF5 "Edited: labor force status" /

TEMIOT "Edited: in the last seven days did you have more than one job?" / TESEMPMNOT "Edited: employment status of spouse or unmarried partner" / TRDPFTPT "Full time or part time employment status of respondent" / TRDTIND1 "Detailed industry record (main job)" / TRTOCC1 "Detailed occupation record (main job)" / TRMIND1 "Intermediate industry record (main job)" / TRMIND1 "Major industry record (main job)" / TRMJOCC1 "Major occupation record (main job)" / TRMJOCCR "Major occupation category (main job)" / TTWK "Weekly earnings topcode flag" / TUOIMPG "Is this business or organization mainly manufacturing, retail trade, wholesale trade, or something else?" (main job)" / TUODP1 "Last time we spoke to someone in this household, you were reported to work for [employer's name]. TUSPPWK "In the last seven days, did your spouse or unmarried partner do any work for pay or profit?" / TUBUS "Does anyone in the household own a business or a farm?" /

```
In [3]: df.columns

Out[3]: Index(['TUCASEID', 'TUFUNENO', 'TUYEAR', 'TUMONTH', 'TEABSR5N', 'TEERN', 'TEERNH10', 'TEERNH2', 'TEERNHRO', 'TEERNHRY', 'TEER', 'TEERNH1', 'TEERNH2', 'TEERNH3', 'TEERNH4', 'TEERNH5', 'TEERNH6', 'TEERNH7', 'TEERNH8', 'TEERNH9', 'TEERNH10', 'TEERNH11', 'TEERNH12', 'TEERNH13', 'TEERNH14', 'TEERNH15', 'TEERNH16', 'TEERNH17', 'TEERNH18', 'TEERNH19', 'TEERNH20', 'TEERNH21', 'TEERNH22', 'TEERNH23', 'TEERNH24', 'TEERNH25', 'TEERNH26', 'TEERNH27', 'TEERNH28', 'TEERNH29', 'TEERNH30', 'TEERNH31', 'TEERNH32', 'TEERNH33', 'TEERNH34', 'TEERNH35', 'TEERNH36', 'TEERNH37', 'TEERNH38', 'TEERNH39', 'TEERNH40', 'TEERNH41', 'TEERNH42', 'TEERNH43', 'TEERNH44', 'TEERNH45', 'TEERNH46', 'TEERNH47', 'TEERNH48', 'TEERNH49', 'TEERNH50', 'TEERNH51', 'TEERNH52', 'TEERNH53', 'TEERNH54', 'TEERNH55', 'TEERNH56', 'TEERNH57', 'TEERNH58', 'TEERNH59', 'TEERNH60', 'TEERNH61', 'TEERNH62', 'TEERNH63', 'TEERNH64', 'TEERNH65', 'TEERNH66', 'TEERNH67', 'TEERNH68', 'TEERNH69', 'TEERNH70', 'TEERNH71', 'TEERNH72', 'TEERNH73', 'TEERNH74', 'TEERNH75', 'TEERNH76', 'TEERNH77', 'TEERNH78', 'TEERNH79', 'TEERNH80', 'TEERNH81', 'TEERNH82', 'TEERNH83', 'TEERNH84', 'TEERNH85', 'TEERNH86', 'TEERNH87', 'TEERNH88', 'TEERNH89', 'TEERNH90', 'TEERNH91', 'TEERNH92', 'TEERNH93', 'TEERNH94', 'TEERNH95', 'TEERNH96', 'TEERNH97', 'TEERNH98', 'TEERNH99', 'TEERNH100', 'TEERNH101', 'TEERNH102', 'TEERNH103', 'TEERNH104', 'TEERNH105', 'TEERNH106', 'TEERNH107', 'TEERNH108', 'TEERNH109', 'TEERNH110', 'TEERNH111', 'TEERNH112', 'TEERNH113', 'TEERNH114', 'TEERNH115', 'TEERNH116', 'TEERNH117', 'TEERNH118', 'TEERNH119', 'TEERNH120', 'TEERNH121', 'TEERNH122', 'TEERNH123', 'TEERNH124', 'TEERNH125', 'TEERNH126', 'TEERNH127', 'TEERNH128', 'TEERNH129', 'TEERNH130', 'TEERNH131', 'TEERNH132', 'TEERNH133', 'TEERNH134', 'TEERNH135', 'TEERNH136', 'TEERNH137', 'TEERNH138', 'TEERNH139', 'TEERNH140', 'TEERNH141', 'TEERNH142', 'TEERNH143', 'TEERNH144', 'TEERNH145', 'TEERNH146', 'TEERNH147', 'TEERNH148', 'TEERNH149', 'TEERNH150', 'TEERNH151', 'TEERNH152', 'TEERNH153', 'TEERNH154', 'TEERNH155', 'TEERNH156', 'TEERNH157', 'TEERNH158', 'TEERNH159', 'TEERNH160', 'TEERNH161', 'TEERNH162', 'TEERNH163', 'TEERNH164', 'TEERNH165', 'TEERNH166', 'TEERNH167', 'TEERNH168', 'TEERNH169', 'TEERNH170', 'TEERNH171', 'TEERNH172', 'TEERNH173', 'TEERNH174', 'TEERNH175', 'TEERNH176', 'TEERNH177', 'TEERNH178', 'TEERNH179', 'TEERNH180', 'TEERNH181', 'TEERNH182', 'TEERNH183', 'TEERNH184', 'TEERNH185', 'TEERNH186', 'TEERNH187', 'TEERNH188', 'TEERNH189', 'TEERNH190', 'TEERNH191', 'TEERNH192', 'TEERNH193', 'TEERNH194', 'TEERNH195', 'TEERNH196', 'TEERNH197', 'TEERNH198', 'TEERNH199', 'TEERNH200', 'TEERNH201', 'TEERNH202', 'TEERNH203', 'TEERNH204', 'TEERNH205', 'TEERNH206', 'TEERNH207', 'TEERNH208', 'TEERNH209', 'TEERNH210', 'TEERNH211', 'TEERNH212', 'TEERNH213', 'TEERNH214', 'TEERNH215', 'TEERNH216', 'TEERNH217', 'TEERNH218', 'TEERNH219', 'TEERNH220', 'TEERNH221', 'TEERNH222', 'TEERNH223', 'TEERNH224', 'TEERNH225', 'TEERNH226', 'TEERNH227', 'TEERNH228', 'TEERNH229', 'TEERNH230', 'TEERNH231', 'TEERNH232', 'TEERNH233', 'TEERNH234', 'TEERNH235', 'TEERNH236', 'TEERNH237', 'TEERNH238', 'TEERNH239', 'TEERNH240', 'TEERNH241', 'TEERNH242', 'TEERNH243', 'TEERNH244', 'TEERNH245', 'TEERNH246', 'TEERNH247', 'TEERNH248', 'TEERNH249', 'TEERNH250', 'TEERNH251', 'TEERNH252', 'TEERNH253', 'TEERNH254', 'TEERNH255', 'TEERNH256', 'TEERNH257', 'TEERNH258', 'TEERNH259', 'TEERNH260', 'TEERNH261', 'TEERNH262', 'TEERNH263', 'TEERNH264', 'TEERNH265', 'TEERNH266', 'TEERNH267', 'TEERNH268', 'TEERNH269', 'TEERNH270', 'TEERNH271', 'TEERNH272', 'TEERNH273', 'TEERNH274', 'TEERNH275', 'TEERNH276', 'TEERNH277', 'TEERNH278', 'TEERNH279', 'TEERNH280', 'TEERNH281', 'TEERNH282', 'TEERNH283', 'TEERNH284', 'TEERNH285', 'TEERNH286', 'TEERNH287', 'TEERNH288', 'TEERNH289', 'TEERNH290', 'TEERNH291', 'TEERNH292', 'TEERNH293', 'TEERNH294', 'TEERNH295', 'TEERNH296', 'TEERNH297', 'TEERNH298', 'TEERNH299', 'TEERNH300', 'TEERNH301', 'TEERNH302', 'TEERNH303', 'TEERNH304', 'TEERNH305', 'TEERNH306', 'TEERNH307', 'TEERNH308', 'TEERNH309', 'TEERNH310', 'TEERNH311', 'TEERNH312', 'TEERNH313', 'TEERNH314', 'TEERNH315', 'TEERNH316', 'TEERNH317', 'TEERNH318', 'TEERNH319', 'TEERNH320', 'TEERNH321', 'TEERNH322', 'TEERNH323', 'TEERNH324', 'TEERNH325', 'TEERNH326', 'TEERNH327', 'TEERNH328', 'TEERNH329', 'TEERNH330', 'TEERNH331', 'TEERNH332', 'TEERNH333', 'TEERNH334', 'TEERNH335', 'TEERNH336', 'TEERNH337', 'TEERNH338', 'TEERNH339', 'TEERNH340', 'TEERNH341', 'TEERNH342', 'TEERNH343', 'TEERNH344', 'TEERNH345', 'TEERNH346', 'TEERNH347', 'TEERNH348', 'TEERNH349', 'TEERNH350', 'TEERNH351', 'TEERNH352', 'TEERNH353', 'TEERNH354', 'TEERNH355', 'TEERNH356', 'TEERNH357', 'TEERNH358', 'TEERNH359', 'TEERNH360', 'TEERNH361', 'TEERNH362', 'TEERNH363', 'TEERNH364', 'TEERNH365', 'TEERNH366', 'TEERNH367', 'TEERNH368', 'TEERNH369', 'TEERNH370', 'TEERNH371', 'TEERNH372', 'TEERNH373', 'TEERNH374', 'TEERNH375', 'TEERNH376', 'TEERNH377', 'TEERNH378', 'TEERNH379', 'TEERNH380', 'TEERNH381', 'TEERNH382', 'TEERNH383', 'TEERNH384', 'TEERNH385', 'TEERNH386', 'TEERNH387', 'TEERNH388', 'TEERNH389', 'TEERNH390', 'TEERNH391', 'TEERNH392', 'TEERNH393', 'TEERNH394', 'TEERNH395', 'TEERNH396', 'TEERNH397', 'TEERNH398', 'TEERNH399', 'TEERNH400', 'TEERNH401', 'TEERNH402', 'TEERNH403', 'TEERNH404', 'TEERNH405', 'TEERNH406', 'TEERNH407', 'TEERNH408', 'TEERNH409', 'TEERNH410', 'TEERNH411', 'TEERNH412', 'TEERNH413', 'TEERNH414', 'TEERNH415', 'TEERNH416', 'TEERNH417', 'TEERNH418', 'TEERNH419', 'TEERNH420', 'TEERNH421', 'TEERNH422', 'TEERNH423', 'TEERNH424', 'TEERNH425', 'TEERNH426', 'TEERNH427', 'TEERNH428', 'TEERNH429', 'TEERNH430', 'TEERNH431', 'TEERNH432', 'TEERNH433', 'TEERNH434', 'TEERNH435', 'TEERNH436', 'TEERNH437', 'TEERNH438', 'TEERNH439', 'TEERNH440', 'TEERNH441', 'TEERNH442', 'TEERNH443', 'TEERNH444', 'TEERNH445', 'TEERNH446', 'TEERNH447', 'TEERNH448', 'TEERNH449', 'TEERNH450', 'TEERNH451', 'TEERNH452', 'TEERNH453', 'TEERNH454', 'TEERNH455', 'TEERNH456', 'TEERNH457', 'TEERNH458', 'TEERNH459', 'TEERNH460', 'TEERNH461', 'TEERNH462', 'TEERNH463', 'TEERNH464', 'TEERNH465', 'TEERNH466', 'TEERNH467', 'TEERNH468', 'TEERNH469', 'TEERNH470', 'TEERNH471', 'TEERNH472', 'TEERNH473', 'TEERNH474', 'TEERNH475', 'TEERNH476', 'TEERNH477', 'TEERNH478', 'TEERNH479', 'TEERNH480', 'TEERNH481', 'TEERNH482', 'TEERNH483', 'TEERNH484', 'TEERNH485', 'TEERNH486', 'TEERNH487', 'TEERNH488', 'TEERNH489', 'TEERNH490', 'TEERNH491', 'TEERNH492', 'TEERNH493', 'TEERNH494', 'TEERNH495', 'TEERNH496', 'TEERNH497', 'TEERNH498', 'TEERNH499', 'TEERNH500', 'TEERNH501', 'TEERNH502', 'TEERNH503', 'TEERNH504', 'TEERNH505', 'TEERNH506', 'TEERNH507', 'TEERNH508', 'TEERNH509', 'TEERNH510', 'TEERNH511', 'TEERNH512', 'TEERNH513', 'TEERNH514', 'TEERNH515', 'TEERNH516', 'TEERNH517', 'TEERNH518', 'TEERNH519', 'TEERNH520', 'TEERNH521', 'TEERNH522', 'TEERNH523', 'TEERNH524', 'TEERNH525', 'TEERNH526', 'TEERNH527', 'TEERNH528', 'TEERNH529', 'TEERNH530', 'TEERNH531', 'TEERNH532', 'TEERNH533', 'TEERNH534', 'TEERNH535', 'TEERNH536', 'TEERNH537', 'TEERNH538', 'TEERNH539', 'TEERNH540', 'TEERNH541', 'TEERNH542', 'TEERNH543', 'TEERNH544', 'TEERNH545', 'TEERNH546', 'TEERNH547', 'TEERNH548', 'TEERNH549', 'TEERNH550', 'TEERNH551', 'TEERNH552', 'TEERNH553', 'TEERNH554', 'TEERNH555', 'TEERNH556', 'TEERNH557', 'TEERNH558', 'TEERNH559', 'TEERNH560', 'TEERNH561', 'TEERNH562', 'TEERNH563', 'TEERNH564', 'TEERNH565', 'TEERNH566', 'TEERNH567', 'TEERNH568', 'TEERNH569', 'TEERNH570', 'TEERNH571', 'TEERNH572', 'TEERNH573', 'TEERNH574', 'TEERNH575', 'TEERNH576', 'TEERNH577', 'TEERNH578', 'TEERNH579', 'TEERNH580', 'TEERNH581', 'TEERNH582', 'TEERNH583', 'TEERNH584', 'TEERNH585', 'TEERNH586', 'TEERNH587', 'TEERNH588', 'TEERNH589', 'TEERNH590', 'TEERNH591', 'TEERNH592', 'TEERNH593', 'TEERNH594', 'TEERNH595', 'TEERNH596', 'TEERNH597', 'TEERNH598', 'TEERNH599', 'TEERNH600', 'TEERNH601', 'TEERNH602', 'TEERNH603', 'TEERNH604', 'TEERNH605', 'TEERNH606', 'TEERNH607', 'TEERNH608', 'TEERNH609', 'TEERNH610', 'TEERNH611', 'TEERNH612', 'TEERNH613', 'TEERNH614', 'TEERNH615', 'TEERNH616', 'TEERNH617', 'TEERNH618', 'TEERNH619', 'TEERNH620', 'TEERNH621', 'TEERNH622', 'TEERNH623', 'TEERNH624', 'TEERNH625', 'TEERNH626', 'TEERNH627', 'TEERNH628', 'TEERNH629', 'TEERNH630', 'TEERNH631', 'TEERNH632', 'TEERNH633', 'TEERNH634', 'TEERNH635', 'TEERNH636', 'TEERNH637', 'TEERNH638', 'TEERNH639', 'TEERNH640', 'TEERNH641', 'TEERNH642', 'TEERNH643', 'TEERNH644', 'TEERNH645', 'TEERNH646', 'TEERNH647', 'TEERNH648', 'TEERNH649', 'TEERNH650', 'TEERNH651', 'TEERNH652', 'TEERNH653', 'TEERNH654', 'TEERNH655', 'TEERNH656', 'TEERNH657', 'TEERNH658', 'TEERNH659', 'TEERNH660', 'TEERNH661', 'TEERNH662', 'TEERNH663', 'TEERNH664', 'TEERNH665', 'TEERNH666', 'TEERNH667', 'TEERNH668', 'TEERNH669', 'TEERNH670', 'TEERNH671', 'TEERNH672', 'TEERNH673', 'TEERNH674', 'TEERNH675', 'TEERNH676', 'TEERNH677', 'TEERNH678', 'TEERNH679', 'TEERNH680', 'TEERNH681', 'TEERNH682', 'TEERNH683', 'TEERNH684', 'TEERNH685', 'TEERNH686', 'TEERNH687', 'TEERNH688', 'TEERNH689', 'TEERNH690', 'TEERNH691', 'TEERNH692', 'TEERNH693', 'TEERNH694', 'TEERNH695', 'TEERNH696', 'TEERNH697', 'TEERNH698', 'TEERNH699', 'TEERNH700', 'TEERNH701', 'TEERNH702', 'TEERNH703', 'TEERNH704', 'TEERNH705', 'TEERNH706', 'TEERNH707', 'TEERNH708', 'TEERNH709', 'TEERNH710', 'TEERNH711', 'TEERNH712', 'TEERNH713', 'TEERNH714', 'TEERNH715', 'TEERNH716', 'TEERNH717', 'TEERNH718', 'TEERNH719', 'TEERNH720', 'TEERNH721', 'TEERNH722', 'TEERNH723', 'TEERNH724', 'TEERNH725', 'TEERNH726', 'TEERNH727', 'TEERNH728', 'TEERNH729', 'TEERNH730', 'TEERNH731', 'TEERNH732', 'TEERNH733', 'TEERNH734', 'TEERNH735', 'TEERNH736', 'TEERNH737', 'TEERNH738', 'TEERNH739', 'TEERNH740', 'TEERNH741', 'TEERNH742', 'TEERNH743', 'TEERNH744', 'TEERNH745', 'TEERNH746', 'TEERNH747', 'TEERNH748', 'TEERNH749', 'TEERNH750', 'TEERNH751', 'TEERNH752', 'TEERNH753', 'TEERNH754', 'TEERNH755', 'TEERNH756', 'TEERNH757', 'TEERNH758', 'TEERNH759', 'TEERNH760', 'TEERNH761', 'TEERNH762', 'TEERNH763', 'TEERNH764', 'TEERNH765', 'TEERNH766', 'TEERNH767', 'TEERNH768', 'TEERNH769', 'TEERNH770', 'TEERNH771', 'TEERNH772', 'TEERNH773', 'TEERNH774', 'TEERNH775', 'TEERNH776', 'TEERNH777', 'TEERNH778', 'TEERNH779', 'TEERNH780', 'TEERNH781', 'TEERNH782', 'TEERNH783', 'TEERNH784', 'TEERNH785', 'TEERNH786', 'TEERNH787', 'TEERNH788', 'TEERNH789', 'TEERNH790', 'TEERNH791', 'TEERNH792', 'TEERNH793', 'TEERNH794', 'TEERNH795', 'TEERNH796', 'TEERNH797', 'TEERNH798', 'TEERNH799', 'TEERNH800', 'TEERNH801', 'TEERNH802', 'TEERNH803', 'TEERNH804', 'TEERNH805', 'TEERNH806', 'TEERNH807', 'TEERNH808', 'TEERNH809', 'TEERNH810', 'TEERNH811', 'TEERNH812', 'TEERNH813', 'TEERNH814', 'TEERNH815', 'TEERNH816', 'TEERNH817', 'TEERNH818', 'TEERNH819', 'TEERNH820', 'TEERNH821', 'TEERNH822', 'TEERNH823', 'TEERNH824', 'TEERNH825', 'TEERNH826', 'TEERNH827', 'TEERNH828', 'TEERNH829', 'TEERNH830', 'TEERNH831', 'TEERNH832', 'TEERNH833', 'TEERNH834', 'TEERNH835', 'TEERNH836', 'TEERNH837', 'TEERNH838', 'TEERNH839', 'TEERNH840', 'TEERNH841', 'TEERNH842', 'TEERNH843', 'TEERNH844', 'TEERNH845', 'TEERNH846', 'TEERNH847', 'TEERNH848', 'TEERNH849', 'TEERNH850', 'TEERNH851', 'TEERNH852', 'TEERNH853', 'TEERNH854', 'TEERNH855', 'TEERNH856', 'TEERNH857', 'TEERNH858', 'TEERNH859', 'TEERNH860', 'TEERNH861', 'TEERNH862', 'TEERNH863', 'TEERNH864', 'TEERNH865', 'TEERNH866', 'TEERNH867', 'TEERNH868', 'TEERNH869', 'TEERNH870', 'TEERNH871', 'TEERNH872', 'TEERNH873', 'TEERNH874', 'TEERNH875', 'TEERNH876', 'TEERNH877', 'TEERNH878', 'TEERNH879', 'TEERNH880', 'TEERNH881', 'TEERNH882', 'TEERNH883', 'TEERNH884', 'TEERNH885', 'TEERNH886', 'TEERNH887', 'TEERNH888', 'TEERNH889', 'TEERNH890', 'TEERNH891', 'TEERNH892', 'TEERNH893', 'TEERNH894', 'TEERNH895', 'TEERNH896', 'TEERNH897', 'TEERNH898', 'TEERNH899', 'TEERNH900', 'TEERNH901', 'TEERNH902', 'TEERNH903', 'TEERNH904', 'TEERNH905', 'TEERNH906', 'TEERNH907', 'TEERNH908', 'TEERNH909', 'TEERNH910', 'TEERNH911', 'TEERNH912', 'TEERNH913', 'TEERNH914', 'TEERNH915', 'TEERNH916', 'TEERNH917', 'TEERNH918', 'TEERNH919', 'TEERNH920', 'TEERNH921', 'TEERNH922', 'TEERNH923', 'TEERNH924', 'TEERNH925', 'TEERNH926', 'TEERNH927', 'TEERNH928', 'TEERNH929', 'TEERNH930', 'TEERNH931', 'TEERNH932', 'TEERNH933', 'TEERNH934', 'TEERNH935', 'TEERNH936', 'TEERNH937', 'TEERNH938', 'TEERNH939', 'TEERNH940', 'TEERNH941', 'TEERNH942', 'TEERNH943', 'TEERNH944', 'TEERNH945', 'TEERNH946', 'TEERNH947', 'TEERNH948', 'TEERNH9
```



## 5.1. Modelling technique

We opted for an unsupervised learning approach to our data, specifically clustering. We chose this method because with unsupervised learning, you can find relationships in data that aren't immediately apparent. It does this by comparing the data based on similarity. In clustering, this is used to group the unlabeled data into clusters that represent patterns that were found in the data.

Our current focus is using K – Mode clustering, an extension of K-Means clustering, to determine the relationships that exist between our variables by presenting them as a group of visually clustered data points which reveal relationships not immediately apparent from the data set. The number of clusters is determined by the K value provided to the method, with several ways to determine the correct K value including an Elbow Chart or a Silhouette Graphs that can provide a visual representation of ideal K values for our data set.

We chose K-Mode clustering because of the nature of our data. K-Mode clusters are extremely efficient when working with large amounts of categorical data. K-Mode doesn't need to calculate all of the pairwise distances between data points. Additionally, it is distribution free, meaning that it does not require imposing distribution assumptions on the data in order to work.

K-Means clustering on the other hand does not work well with categorical data sets, primarily because it has the issues above that K-Mode does not, and ultimately it's objective function simply doesn't work with categorical data.

```
In [ ]: # First we will keep a copy of data
df_employment_copy = df_employment2.copy()
# Doing k-mode with "kmode" initialization
km_cao = KModes(n_clusters=2, init = "Cao", n_init = 1, verbose=1)
fitClusters_cao = km_cao.fit_predict(df_employment2)

In [ ]: # Predicted Clusters
fitClusters_cao

In [ ]: clusterCentroidsDf = pd.DataFrame(km_cao.cluster_centroids_)
clusterCentroidsDf.columns = df_employment2.columns

In [ ]: # Mode of the clusters
clusterCentroidsDf

In [ ]: # Using K-Mode with "huang" initialization
km_huang = KModes(n_clusters=2, init = "huang", n_init = 1, verbose=1)
fitClusters_huang = km_huang.fit_predict(df_employment2)

In [ ]: # Predicted Clusters
fitClusters_huang
```

## Building the clusters using K-Mode clustering.

- Clustering is a means of grouping data based on characteristics. The clusters will form based on relationships between the individual data points and potentially reveal new relationships that were not previously apparent.
- Clustering also allows for data compression when working with large data sets. After the data is clustered it can be referred to by cluster ID rather than the individual data points themselves.
- Another use of clustering is the preserve the privacy of the respondent. Data grouped in clusters is practically anonymized by its presence in the cluster and future reference by cluster ID. In our example, no individual respondent ID will need to be used going forward as we can refer to the groups by their cluster IDs.
- Here we used K-Mode clustering, which is similar to K - Means clustering but is based on the Mode of K. K-Mode clusters are a good choice when working with large amounts of categorical data like what is present in our data set.

```
In [ ]: # Choosing K by comparing Cost against each K
cost = []
for num_clusters in list(range(1,10)):
    kmode = KModes(n_clusters=num_clusters, init = "huang", n_init = 1, verbose=1)
    kmode.fit_predict(df_employment2)
    cost.append(kmode.cost_)

In [ ]: #For KModes, plot cost for a range of K values. Cost is the sum of all the dissimilarities between the clusters.
# Select the K where you observe an elbow-like bend with a lesser cost value.
plt.subplots(figsize = (15,5))
y = np.array([i for i in range(1,10,1)])
plt.plot(y, cost)
```

- Here we used what is called the Elbow Method for determining the optimal value of K which will also determine the number of clusters for our purposes. When examining the chart visually, one looks for the point at which the line makes an angle similar to that of an arm bent at the elbow. In our case there are a few bends, but the one with the least dispersion associated with it is at K=3, and as such that was the value we went forward with.

```
In [ ]: ## Choosing K=3
km_cao = KModes(n_clusters=3, init = "Cao", n_init = 1, verbose=1)
fitClusters_cao = km_cao.fit_predict(df_employment2)

In [ ]: fitClusters_cao.shape

In [ ]: # Combining the predicted clusters with the original DF
df_employment2 = df_employment2_copy.reset_index()

In [ ]: clustersDf = pd.DataFrame(fitClusters_cao)
clustersDf.columns = ['cluster_predicted']
combinedDf = pd.concat([df_employment2, clustersDf], axis = 1).reset_index()
combinedDf = combinedDf.drop(['index', 'level_0'], axis = 1)

In [ ]: combinedDf.head()

In [ ]: # Cluster Identification
cluster_0 = combinedDf[combinedDf['cluster_predicted'] == 0]
cluster_1 = combinedDf[combinedDf['cluster_predicted'] == 1]
cluster_2 = combinedDf[combinedDf['cluster_predicted'] == 2]

In [ ]: ## Plotting the predicted cluster for first feature only
'''
i = 'TEENMFR_1_0'
plt.subplots(figsize = (15,5))
sns.countplot(x=combinedDf[i],order=combinedDf[i].value_counts().index,hue=combinedDf['cluster_predicted'])
plt.tight_layout()
plt.show()

## Plotting the predicted cluster for each feature
for i in df_employment3.columns.values:
    if i != 'index':
        continue
    else:
        plt.subplots(figsize = (15,5))
        sns.countplot(x=combinedDf[i],order=combinedDf[i].value_counts().index,hue=combinedDf['cluster_predicted'])
        plt.tight_layout()
        plt.show()
# comment out break statement to generate plots for all features
break
```

- As the above plot shows, per predicted feature the survey respondent with hourly earnings are grouped in cluster 0 while cluster 1 contains the mix of respondent with an hourly earning & non hourly earnings.

## 5.4. Assess Model

The silhouette Method is also a method to find the optimal number of clusters and interpretation and validation of consistency within clusters of data. The silhouette method computes silhouette coefficients of each point that measure how much a point is similar to its own cluster compared to other clusters, by providing a succinct graphical representation of how well each object has been classified.

**Model assessment** - Summarise the results of this task, list the qualities of your generated models (e.g.in terms of accuracy) and rank their quality in relation to each other.

**Revised parameter settings** - According to the model assessment, revise parameter settings and tune them for the next modelling run. Iterate model building and assessment until you strongly believe that you have found the best model(s). Document all such revisions and assessments.

```
In [ ]: X = df_employment2
from sklearn.cluster import KMeans
from sklearn.metrics import silhouette_samples, silhouette_score

import matplotlib.pyplot as plt
import matplotlib.cm as cm

n_clusters = 10
range_n_clusters = [2, 3, 4, 5, 6, 7, 8]
for n_clusters in range_n_clusters:
    # Create a subplot with 1 row and 2 columns
    fig, (ax1, ax2) = plt.subplots(1, 2)
    fig.set_size_inches(18, 7)

    # The 1st subplot is the silhouette plot
    # The silhouette coefficient can range from -1, 1 but in this example all
    # lie within [-0.5, 1]
    ax1.set_xlim((-0.1, 1))
    # The (n_clusters+1)*10 is for inserting blank space between silhouette
    # plots of individual clusters, to demarcate them clearly.
    ax1.set_ylim([0, len(X) * (n_clusters + 1) * 10])

    # Initialize the clusterer with n_clusters value and a random generator
    # seed of 10 for reproducibility.
    clusterer = KMeans(n_clusters=n_clusters, random_state=10)
    cluster_labels = clusterer.fit_predict(X)

    # The silhouette score gives the average value for all the samples.
    # This gives a perspective into the density and separation of the formed
    # clusters
    silhouette_avg = silhouette_score(X, cluster_labels)
    print("For n_clusters =", n_clusters,
          "The average silhouette score is :", silhouette_avg)

    # Compute the silhouette scores for each sample
    sample_silhouette_values = silhouette_samples(X, cluster_labels)

    y_lower = 10
    for i in range(n_clusters):
        # Aggregate the silhouette scores for samples belonging to
        # cluster i, and sort them
        ith_cluster_silhouette_values = \
            sample_silhouette_values[cluster_labels == i]

        ith_cluster_silhouette_values.sort()

        size_cluster_i = ith_cluster_silhouette_values.shape[0]
        y_upper = y_lower + size_cluster_i

        color = cm.nipy_spectral(float(i) / n_clusters)
        ax1.fill_between(plt.arange(y_lower, y_upper),
            0, ith_cluster_silhouette_values,
            facecolor=color, edgecolor=color, alpha=0.7)

        # Label the silhouette plots with their cluster numbers at the middle
        ax1.text(-0.05, y_lower + 0.5 * size_cluster_i, str(i))

        # Compute the new y lower for next plot
        y_lower = y_upper + 10 # 10 for the 0 samples

    ax1.set_title("The silhouette plot for the various clusters.")
    ax1.set_xlabel("The silhouette coefficient values")
    ax1.set_ylabel("Cluster label")

    # The vertical line for average silhouette score of all the values
    ax1.axvline(x=silhouette_avg, color="red", linestyle="--")

    ax1.set_yticks([]) # Clear the yaxis labels & ticks
    ax1.set_xticks((-0.1, 0, 0.2, 0.4, 0.6, 0.8, 1))
    plt.show()
```

- Silhouette analysis can be used to study the separation distance between clusters. The silhouette plot displays a measure of how close each point in one cluster is to points in the neighboring clusters and thus provides a way to assess parameters like number of clusters.
- Observations from above Silhouette Plots:
- The silhouette plot shows that the n\_cluster value of 3 is a bad pick, as all the points in the cluster with cluster\_label=0 are below-average silhouette scores.
- The silhouette plot shows that the n\_cluster value of 5 is a bad pick, as all the points in the cluster with cluster\_label=2 and 4 are below-average silhouette scores.
- The silhouette plot shows that the n\_cluster value of 6 is a bad pick, as all the points in the cluster with cluster\_label=1,2,4 and 5 are below-average silhouette scores, and also due to the presence of outliers.
- Silhouette analysis is more ambivalent in deciding between 2 and 4.
- The thickness of the silhouette plot for the cluster with cluster\_label=1 when n\_clusters=2, is bigger in size owing to the grouping of the 3 sub-clusters into one big cluster.
- For n\_clusters=4, all the plots are more or less of similar thickness and hence are of similar sizes, as can be considered as best 'K'.

```
In [ ]: #extract features and target variables
x = combinedDf.drop(columns='cluster_predicted')
y = combinedDf['cluster_predicted']
#save the feature name and target variables
feature_names = x.columns
labels = y.unique()
#split the dataset
from sklearn.model_selection import train_test_split
X_train, test_x, y_train, test_lab = train_test_split(X,y,
                                                    test_size = 0.4,
                                                    random_state = 42)
```

- In the above cell, the data is being divided into multiple sub sets in order to both train and then subsequently test the our modeling solution against. This separation into multiple sets is an in effort to avoid the possibility of over fitting our model against the training set and thus skewing our results.

```
In [ ]: from sklearn.tree import DecisionTreeClassifier
clf = DecisionTreeClassifier(max_depth=3, random_state = 42)
clf.fit(X_train, y_train)

In [ ]: #import relevant functions
from sklearn.tree import export_text
#export the decision tree
tree_rules = export_text(clf,
                          feature_names = list(feature_names))
#print the result
print(tree_rules)
```

- In the above cell we've created a decision tree to help determine which of our variables are associated with our other variables and how. Decision trees work sort of like flow charts of possible outcomes based on conditions that all start out at one root node. In this case the root of tree is the variable "TEIO1CDD" and all of the branches of our tree extend from it.

```
In [ ]: test_pred_decision_tree = clf.predict(test_x)

In [ ]: #import the relevant packages
from sklearn import metrics
import seaborn as sns
import matplotlib.pyplot as plt
#get the confusion matrix
confusion_matrix = metrics.confusion_matrix(test_lab,
                                             test_pred_decision_tree)

#turn this into a dataframe
matrix_df = pd.DataFrame(confusion_matrix)
#plot the result
ax = plt.axes()
sns.set(font_scale=1.3)
plt.figure(figsize=(10,7))
#set axis titles
ax.set_title('Confusion Matrix - Decision Tree')
ax.set_xlabel('Predicted labels', fontsize=15)
#set xtick labels([''])+labels)
ax.set_xticklabels(['True label'], rotation=15)
ax.set_ylabel('True label', fontsize=15)
ax.set_yticklabels(list(labels), rotation = 0)
plt.show()
```

- The above visual is called a confusion matrix. These closely resemble heatmaps and are used predominantly to help with statistical classification.
- Confusion matrices provide a way to visually observe the accuracy of a classification model based built from a know ground truth and branching into predicted values.

## Assessing the accuracy of our completed model

```
In [ ]: metrics.accuracy_score(test_lab, test_pred_decision_tree)

In [ ]: print(metrics.classification_report(test_lab,
                                       test_pred_decision_tree))

In [ ]: #extract importance
importance = pd.DataFrame({'feature': X_train.columns,
                          'importance': np.round(metrics.feature_importances_, 3)})
importance.sort_values('importance', ascending=False, inplace = True)
print(importance)

In [ ]: from sklearn.model_selection import GridSearchCV
tuned_parameters = [{'max_depth': [1,2,3,4,5],
                     'min_samples_split': [2,4,6,8,10]}]
scores = {'recall':[]}
for score in scores:
    print()
    print(f"Finding hyperparameters for {score}")
    print()
    clf = GridSearchCV(
        DecisionTreeClassifier(), tuned_parameters,
        scoring = f'({score})_macro')
    clf.fit(X_train, y_train)

    print("Best parameters set found on development set:")
    print()
    print(clf.best_params_)
    print()
    print("Grid scores on development set:")
    means = clf.cv_results_['mean_test_score']
    stds = clf.cv_results_['std_test_score']
    for mean, std, params in zip(means, stds,
                                print(f"mean{0.3f} (+/-std{2:0.3f}) for {params}")
```

- The above indicates that our predictive model, as currently constructed has an accuracy of .97 or 97%. This indicates that our model is an excellent predictor of outcomes based on the relationships between our variables and which of them is likely to be indicative of another.

## Model Ensemble

```
In [ ]: ## Train Model
# Import the model we are using
from sklearn.ensemble import RandomForestClassifier
# Instantiate model with 10 decision trees
rf = RandomForestClassifier(n_estimators = 10, random_state = 42)
# Train the model on training data
rf.fit(X_train, y_train)

In [ ]: # Use the forest's predict method on the test data
predictions = rf.predict(test_x)

In [ ]: #import scikit-learn metrics module for accuracy calculation
from sklearn import metrics
# Model Accuracy, how often is the classifier correct?
print("Accuracy:",metrics.accuracy_score(test_lab, predictions))

In [ ]: #### Finding Important Features in Skikit-learn

In [ ]: feature_list = list(X_train.columns)
feature_imp = pd.Series(rf.feature_importances_, index=feature_list).sort_values(ascending=False)
feature_imp_top20 = feature_imp.head(20)
feature_imp.head(20)

In [ ]: import matplotlib.pyplot as plt
import seaborn as sns
import matplotlib inline
# Creating a bar plot
plt.figure(figsize=(10,7))
sns.barplot(x=feature_imp_top20, y=feature_imp_top20.index)
# Add labels to your graph
plt.xlabel('Feature Importance Score')
plt.ylabel('Features')
plt.title("Visualizing Important Features")
plt.legend()
plt.show()
```

## Conclusion

Members of a sub class of the data set are best fits for targeting for advertisement of financial services products based on demographic data revealed in the survey used. The members of that sub class are easily identified by our algorithm for the purposes of identifying a target candidate pool to present to our customers. This suitability is based on a number of key factors such as employment and income level which indicate potential need or desire for specific financial services such as banks accounts or Certificates of Deposit (CDs).

## Next Steps

- Business Application:

Now that we've built and assessed the accuracy of our model the next step is to apply it and interpret the results of our findings and how they apply to our ability to answer our questions. Based on the intermediate results of our model, it is clear that the clusters we have created are indicative of the predictive value of our variables. Going forward, we will determine the profile of the ideal candidate based on the output of our predictive modeling, identifying specifically which traits are most indicative of a candidate's suitability for our services and be able to provide that information to clients for use in their marketing campaigns.

## References

- Data Science in 5 minutes: What is one hot encoding? Educative. (n.d.). Retrieved February 9, 2022, from <https://www.educative.io/blog/one-hot-encoding>
- Decision trees in python - step-by-step implementation. AskPython. (2020, December 7). Retrieved February 10, 2022, from <https://www.askpython.com/python/examples/decision-trees>
- Franklin, S. J. (2019, November 26). Elbow method of k-means clustering algorithm. Medium. Retrieved February 9, 2022, from <https://medium.com/analytixs-vdhy/elbow-method-of-k-means-clustering-algorith-a0c916adc540>
- Google. (n.d.). Clustering algorithms | clustering in machine learning | google developers. Google. Retrieved February 10, 2022, from <https://developers.google.com/machine-learning/clustering/clustering-algorithms>
- Jaiswal, S. (2020, July 10). K-modelsclustering. Medium. Retrieved February 10, 2022, from <https://medium.com/@shailjanitp2013/k-modelsclustering-ef69ef06449>
- Mean imputation for missing data (example in R & SPSS). Statistics Globe. (2022, January 19). Retrieved February 9, 2022, from <https://statisticsglobe.com/mean-imputation-for-missing-data/>
- Selecting the number of clusters with silhouette analysis on kmeans clustering. scikit. (n.d.). Retrieved February 9, 2022, from [https://scikit-learn.org/stable/auto\\_examples/cluster/plot\\_kmeans\\_silhouette\\_analysis.html](https://scikit-learn.org/stable/auto_examples/cluster/plot_kmeans_silhouette_analysis.html)
- U.S. Bureau of Labor Statistics. (n.d.). Atus News releases. U.S. Bureau of Labor Statistics. Retrieved February 10, 2022, from <https://www.bls.gov/tus/>

```
In [ ]:
```