# HONEYBEE PESTS & PATHOGENS IN ONTARIO APIARIES

Moganaviniith Rathinavel

Ragavi Mudaliyar

Paras Gangani

# Table of Contents

# About project

To create a report and do a predictive analysis on pests and pathogens level in apiaries of honeybee in a particular province, Ontario. The prevalence and load (levels or intensity) of pathogens at various times during the beekeeping season was assessed.
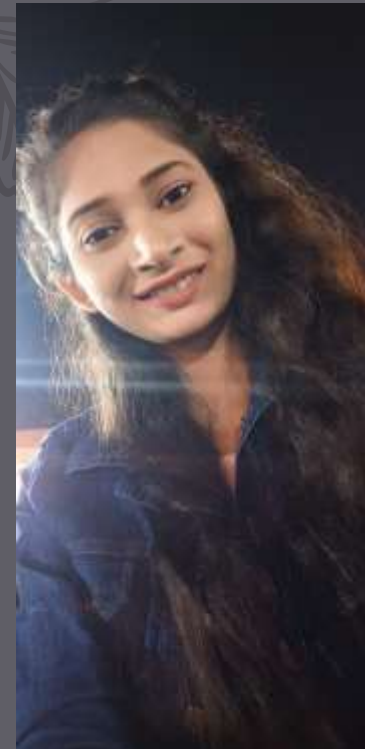
# Meet our Team



**Rajalakshmi Nagarajan**

Team lead/Developer

**Ragavi Mudaliyar**

Communicator
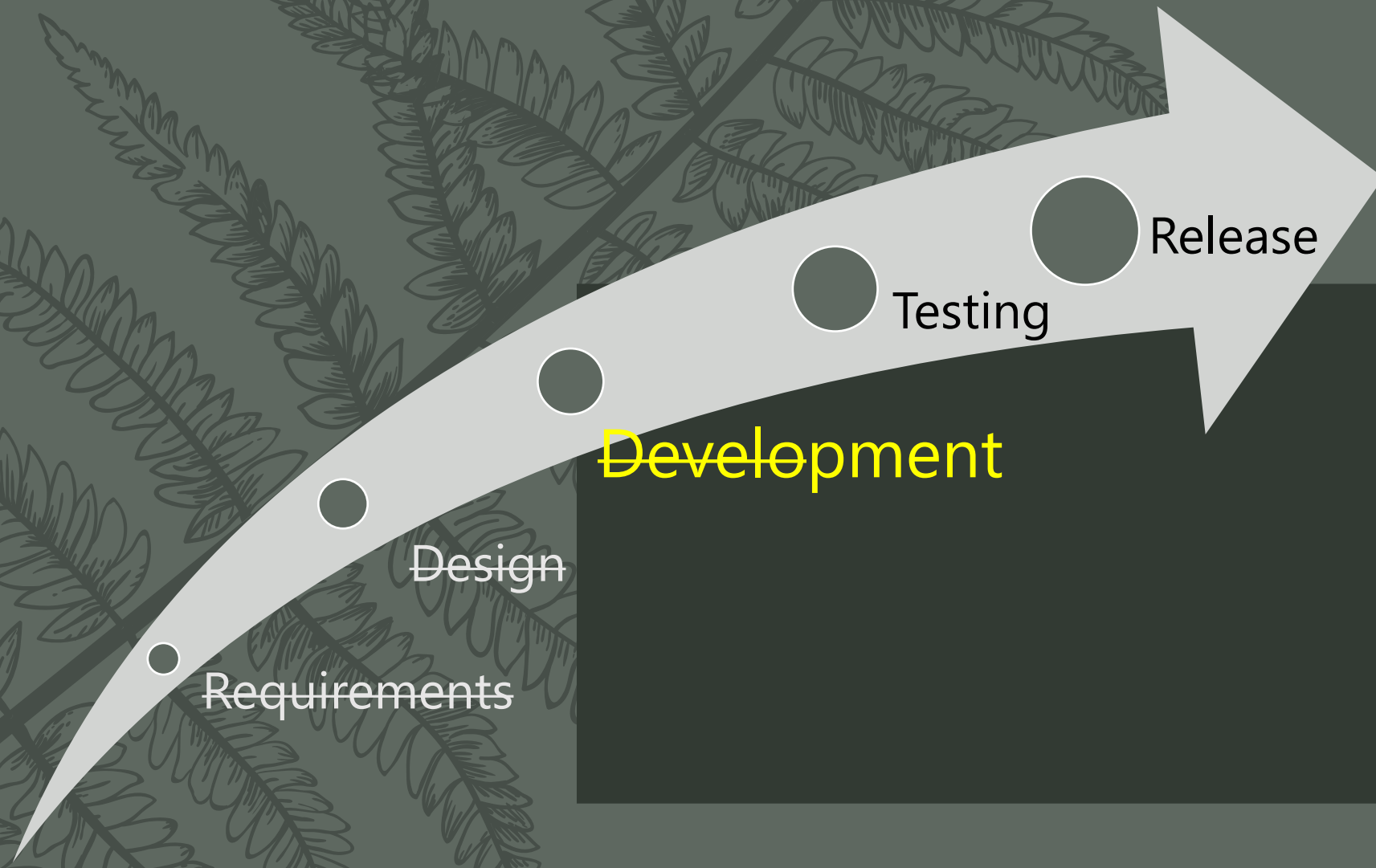
**Paras Gangani**

Analyst

# Project Status



## On Track

**Date of report: 20th February. 2023**

**Date of last report: 15th February. 2023**

# Project Phase

Requirements

Design

Development

Testing

Release

# Milestones

**Jan - Feb**

Acceptance of Project charter — Jan'04

Requirements Phase — Jan'11

Design Phase — Jan'18

Development Phase-I — Feb'08

Development Phase-II — Feb'15

Testing Phase - I — Feb'22

**Mar - Apr**

Testing Phase - II — Mar'01

Release Phase — Mar'08

Client Presentation — Mar'15

Handover of the Deliverables Phase — Mar'22

Project Close-Out Phase — Mar'29

Deliver to client — Apr'5

Apr'12

# ① Requirements

# About the Dataset

Website reference: https://data.ontario.ca/en/dataset/honey-bee-pests-and-pathogens-in-ontario-apiaries

# ② Design

# Documentation

**1** Project Proposal     Microsoft Word Document

**2** Project Charter     Microsoft Word Document

**3** Sharepoint Site     https://georgiancollege.sharepoint.com/sites/HoneybeepestsandpathogensinOntario

③ Development

# Datatypes of variables and missing values distribution for year 2019

```
# check datatype in each column
print("Column datatypes: ")
print(honeybee_2019.dtypes)
```

```
Column datatypes:
Monitoring Site                                                   int64
Inspection Period                                                 int64
Inspection Start Date                                             object
Collection Date                                                   object
Region                                                            object
County                                                            object
Num. Colonies Inspected                                          float64
Num. Colonies - No AFB Found                                     float64
Num. Colonies with AFB (< 10 Cells)                             float64
Num. Colonies with AFB (10 or More Cells)                       float64
Num. Colonies - No EFB Found                                     float64
Num. Colonies with EFB (< 10 Cells)                             float64
Num. Colonies with EFB (10 or More Cells)                       float64
Num. Colonies - No Chalkbrood Found                             float64
Num. Colonies with Chalkbrood (< 10 Cells)                      float64
Num. Colonies with Chalkbrood (10 or More Cells)                float64
Num. Colonies - No Sacbrood Found                               float64
Num. Colonies with Sacbrood (< 10 Cells)                        float64
Num. Colonies with Sacbrood (10 or More Cells)                  float64
Num. Colonies with SHB Adults (1-20)                            float64
Num. Colonies with SHB Adults (>20)                             float64
Num. Colonies with SHB Larvae (1-20)                            float64
Num. of Colonies with SHB Larvae (21-1/4cup)                    float64
Num. Colonies with SHB Larvae (>1/4 cup)                        float64
Average Varroa Infestation (%)                                  float64
Max Varroa Infestation (%)                                      float64
Num. Colonies - Queenless                                       float64
Num. Colonies - Queenright                                      float64
Num. Colonies - Queen Newly Installed                           float64
Num. Colonies - Virgin Queen                                    float64
Num. Colonies - Queen Not Observed                              float64
% Colonies Queenless in Yard at Inspection                       object
Acute Bee Paralysis Virus (log10 RNA copies/bee) - Average      float64
Deformed Wing Virus (log10 RNA copies/bee) - Average            float64
Israeli Acute Paralysis Virus (log10 RNA copies/bee) - Average  float64
Nosema ceranae (log10 DNA copies/bee) - Average                 float64
Kashmir Bee Virus (log10 RNA copies/bee)                        float64
Sacbrood Virus (log10 RNA copies/bee)                           float64
Tracheal Mite Infestation (# bees infested per 25 bees tested)   int64
dtype: object
```

```
# examining missing values
print("Missing values distribution: ")
print(honeybee_2019.isnull().mean())
print("")
```

```
Missing values distribution:
Monitoring Site                                                 0.000000
Inspection Period                                               0.000000
Inspection Start Date                                           0.010989
Collection Date                                                 0.000000
Region                                                          0.000000
County                                                          0.000000
Num. Colonies Inspected                                         0.010989
Num. Colonies - No AFB Found                                    0.010989
Num. Colonies with AFB (< 10 Cells)                            1.000000
Num. Colonies with AFB (10 or More Cells)                      1.000000
Num. Colonies - No EFB Found                                   0.010989
Num. Colonies with EFB (< 10 Cells)                            0.989011
Num. Colonies with EFB (10 or More Cells)                      1.000000
Num. Colonies - No Chalkbrood Found                            0.010989
Num. Colonies with Chalkbrood (< 10 Cells)                     0.901099
Num. Colonies with Chalkbrood (10 or More Cells)               0.802198
Num. Colonies - No Sacbrood Found                              0.010989
Num. Colonies with Sacbrood (< 10 Cells)                       0.989011
Num. Colonies with Sacbrood (10 or More Cells)                 0.989011
Num. Colonies with SHB Adults (1-20)                           1.000000
Num. Colonies with SHB Adults (>20)                            1.000000
Num. Colonies with SHB Larvae (1-20)                           1.000000
Num. of Colonies with SHB Larvae (21-1/4cup)                   1.000000
Num. Colonies with SHB Larvae (>1/4 cup)                       1.000000
Average Varroa Infestation (%)                                 0.010989
Max Varroa Infestation (%)                                     0.010989
Num. Colonies - Queenless                                      0.813187
Num. Colonies - Queenright                                     0.010989
Num. Colonies - Queen Newly Installed                          0.934066
Num. Colonies - Virgin Queen                                   0.945055
Num. Colonies - Queen Not Observed                             1.000000
% Colonies Queenless in Yard at Inspection                     0.010989
Acute Bee Paralysis Virus (log10 RNA copies/bee) - Average     0.000000
Deformed Wing Virus (log10 RNA copies/bee) - Average           0.000000
Israeli Acute Paralysis Virus (log10 RNA copies/bee) - Average 0.000000
Nosema ceranae (log10 DNA copies/bee) - Average                0.000000
Kashmir Bee Virus (log10 RNA copies/bee)                       0.000000
Sacbrood Virus (log10 RNA copies/bee)                          0.000000
Tracheal Mite Infestation (# bees infested per 25 bees tested) 0.000000
dtype: float64
```

# Cleaning Dataset - 2019

```python
# cleaning the column outliers
columns = ['Num. Colonies with Chalkbrood (< 10 Cells)', 'Num. Colonies with Chalkbrood (10 or More Cells)',
           'Num. Colonies - Queenless', 'Num. Colonies - Queen Newly Installed', 'Num. Colonies - Virgin Queen']

# Looping through the columns to fill the entries with NaN values with 0
for column in columns:
    df[column] = df[column].fillna(0)
```

```python
# Convert the dictionary into DataFrame
df = pd.DataFrame(honeybee_2019)
# Remove columns with no values
df = df.drop(['Num. Colonies with AFB (< 10 Cells)', 'Num. Colonies with AFB (10 or More Cells)',
              'Num. Colonies with EFB (< 10 Cells)', 'Num. Colonies with EFB (10 or More Cells)',
              'Num. Colonies with Sacbrood (< 10 Cells)', 'Num. Colonies with Sacbrood (10 or More Cells)',
              'Num. Colonies with SHB Adults (1-20)', 'Num. Colonies with SHB Adults (>20)', 'Num. Colonies with SHB Larvae (1-20)',
              'Num. of Colonies with SHB Larvae (21-1/4cup)', 'Num. Colonies with SHB Larvae (>1/4 cup)',
              'Num. Colonies - Queen Not Observed'], axis=1)
```
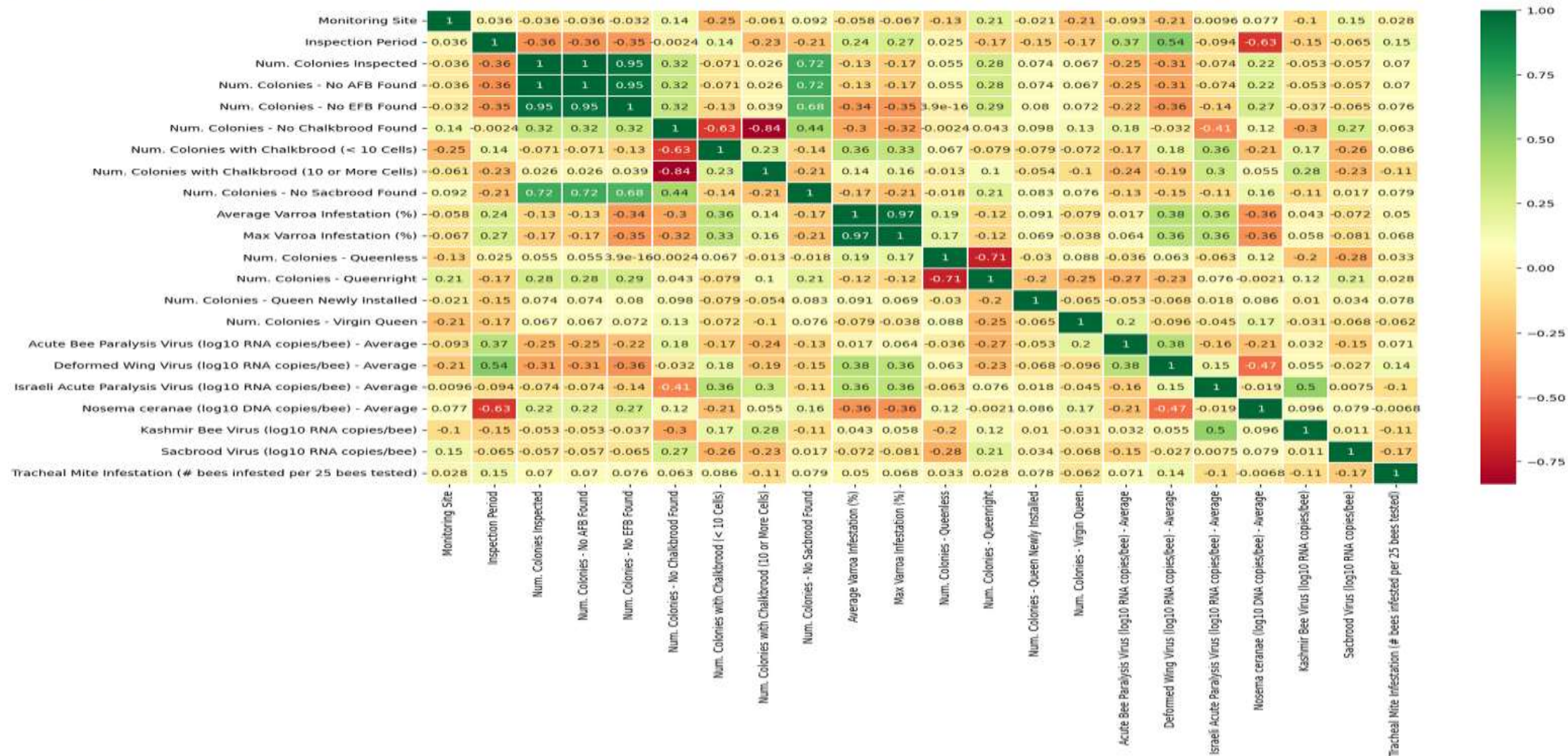
```python
df.head()
```

| | Monitoring Site | Inspection Period | Inspection Start Date | Collection Date | Region | County | Num. Colonies Inspected | Num. Colonies - No AFB Found | Num. Colonies - No EFB Found | Num. Colonies - No Chalkbrood Found | ... | Num. Colonies - Queen Newly Installed | Num. Colonies - Virgin Queen | % Colonies Queenless in Yard at Inspection | Acute Paraly V (log10 F copies/t - Aver |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 06-27-19 | 2019-06-27 | East | LENNOX & ADDINGTON COUNTY | 6.0 | 6.0 | 6.0 | 3.0 | ... | 0.0 | 0.0 | 0% | 0. |
| 1 | 1 | 2 | 08-29-19 | 2019-08-29 | East | LENNOX & ADDINGTON COUNTY | 6.0 | 6.0 | 6.0 | 1.0 | ... | 0.0 | 0.0 | 16.7% | 0. |
| 2 | 1 | 3 | 09-24-19 | 2019-09-24 | East | LENNOX & ADDINGTON COUNTY | 6.0 | 6.0 | 6.0 | 3.0 | ... | 0.0 | 0.0 | 0% | 0. |
| 3 | 2 | 1 | 06-11-19 | 2019-06-11 | South | HALTON REGION | 6.0 | 6.0 | 6.0 | 6.0 | ... | 0.0 | 0.0 | 0% | 0. |
| 4 | 2 | 2 | 08-12-19 | 2019-08-12 | South | HALTON REGION | 6.0 | 6.0 | 6.0 | 6.0 | ... | 0.0 | 0.0 | 0% | 6. |

5 rows × 27 columns

```
In [53]: df['County'].value_counts()

Out[53]: SIMCOE COUNTY                             12
         MIDDLESEX COUNTY                           6
         LAMBTON                                    6
         OXFORD COUNTY                              6
         GREY COUNTY                                6
         FRONTENAC                                  6
         WATERLOO REGION                            6
         LENNOX & ADDINGTON COUNTY                  3
         WELLINGTON                                 3
         PEEL REGION                                3
         STORMONT, DUNDAS & GLENGARRY COUNTY        3
         OTTAWA REGION                              3
         ELGIN COUNTY                               3
         HAMILTON REGION                            3
         NORFOLK COUNTY                             3
         NIPISSING DISTRICT                         3
         DURHAM REGION                              3
         LEEDS & GRENVILLE COUNTY                   3
         THUNDER BAY DISTRICT                       2
         HURON COUNTY                               2
         HALTON REGION                              2
         BRUCE COUNTY                               2
         YORK REGION                                1
         Name: County, dtype: int64
```
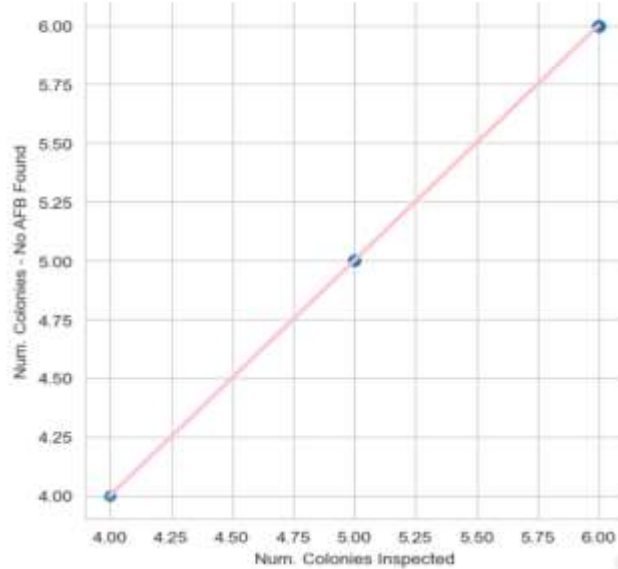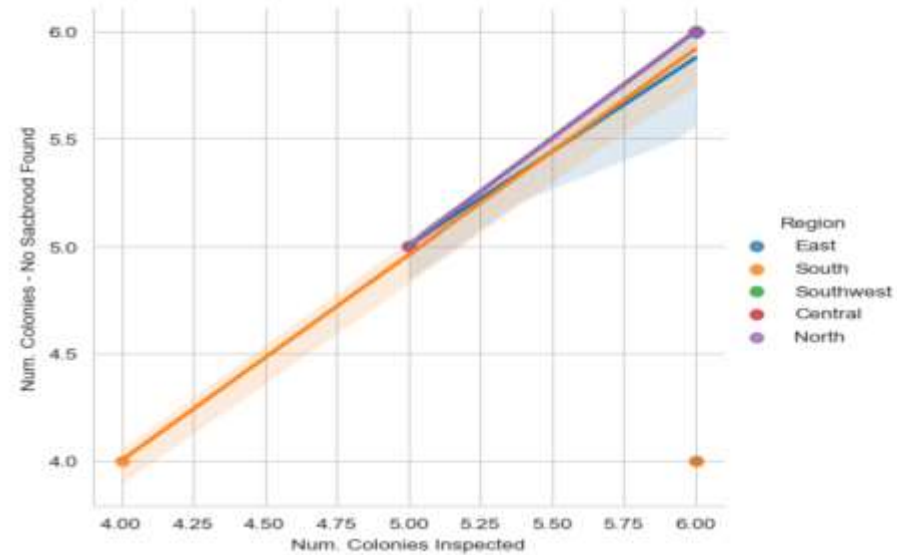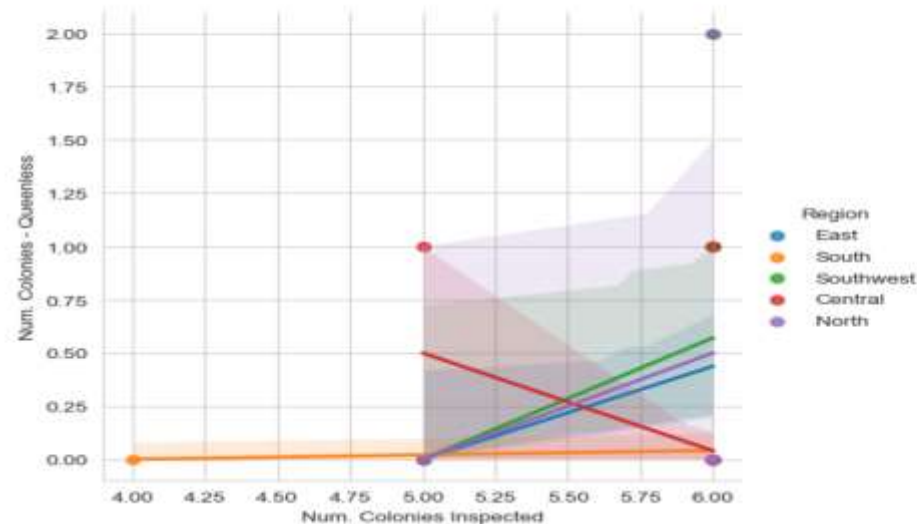
16

# Predictive analysis - K means clustering
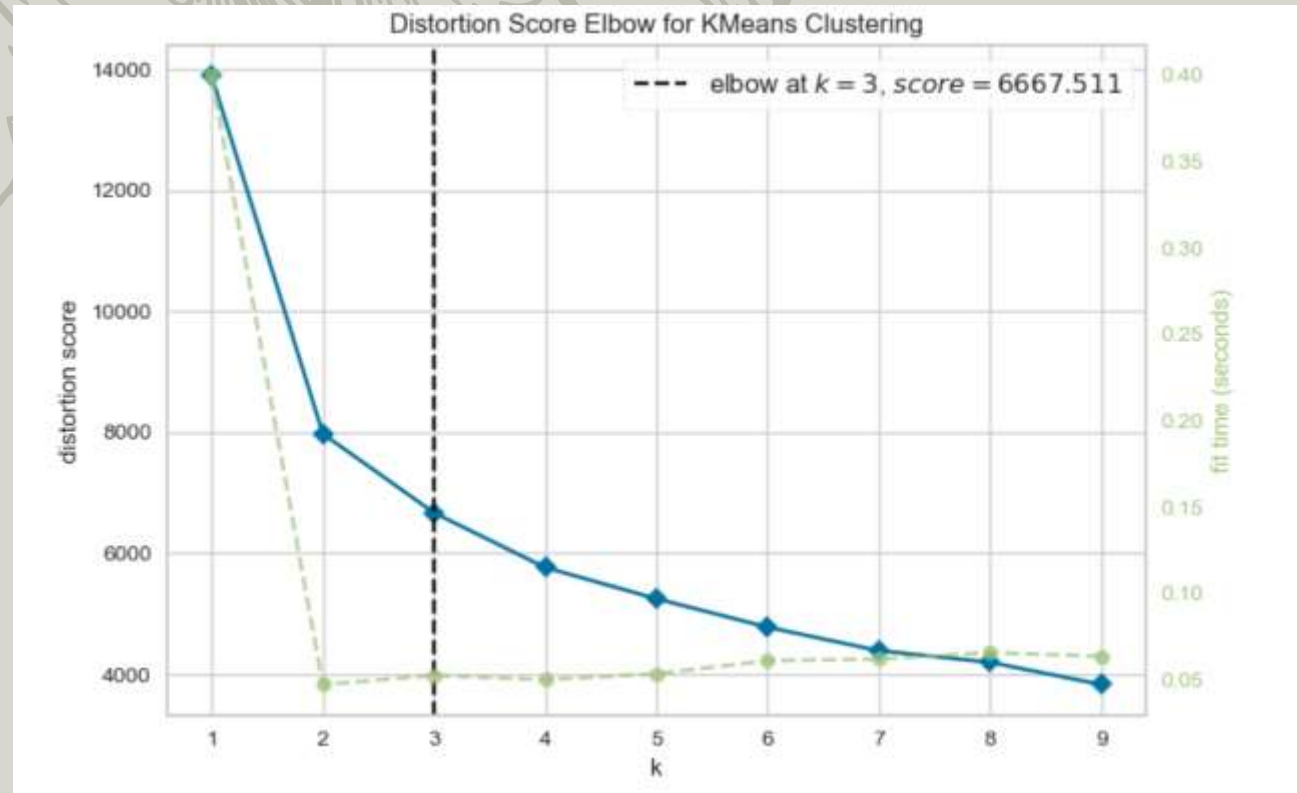
```python
import plotly.express as px

from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler, Normalizer
from sklearn.cluster import KMeans
from sklearn import metrics

scaler = StandardScaler()
# scaler = Normalizer()
df_copy_scaled = scaler.fit_transform(df_copy)

pca = PCA(2, random_state=42)
df_copy_pca = pca.fit_transform(df_copy_scaled)

projection = pd.DataFrame(columns=['x','y'], data=df_copy_pca)
projection
```

|    | x         | y         |
|----|-----------|-----------|
| 0  | 0.456382  | 4.092742  |
| 1  | 4.241694  | 5.907613  |
| 2  | 2.485419  | 3.643503  |
| 3  | -1.773283 | -0.172087 |
| 4  | -1.059917 | -1.009002 |
| 5  | -0.562994 | -0.578608 |
| 6  | 0.369572  | -1.631928 |
| 7  | 1.112740  | -3.051635 |
| 8  | -1.323935 | -0.316031 |
| 9  | -0.287655 | 0.028247  |
| 10 | 0.289717  | -0.497782 |



Distortion Score Elbow for KMeans Clustering

--- elbow at $k = 3$, $score = 6667.511$

# Predictive analysis - K means clustering

# Project Cost

## Issues or Challenges encountered this week and what was done to overcome them

We are using Microsoft Excel for cleaning and grouping of data.

**Update(25'jan):** We are using Python for data cleaning instead of doing manually in Excel.

Understanding outliers and cleaning the data is quite challenging.

Data of years 2017, 2018 and 2019 are considered.

**Update(01'feb):** no challenges

**Update(08'feb):** Understanding the numerical data visualization is quite challenging.
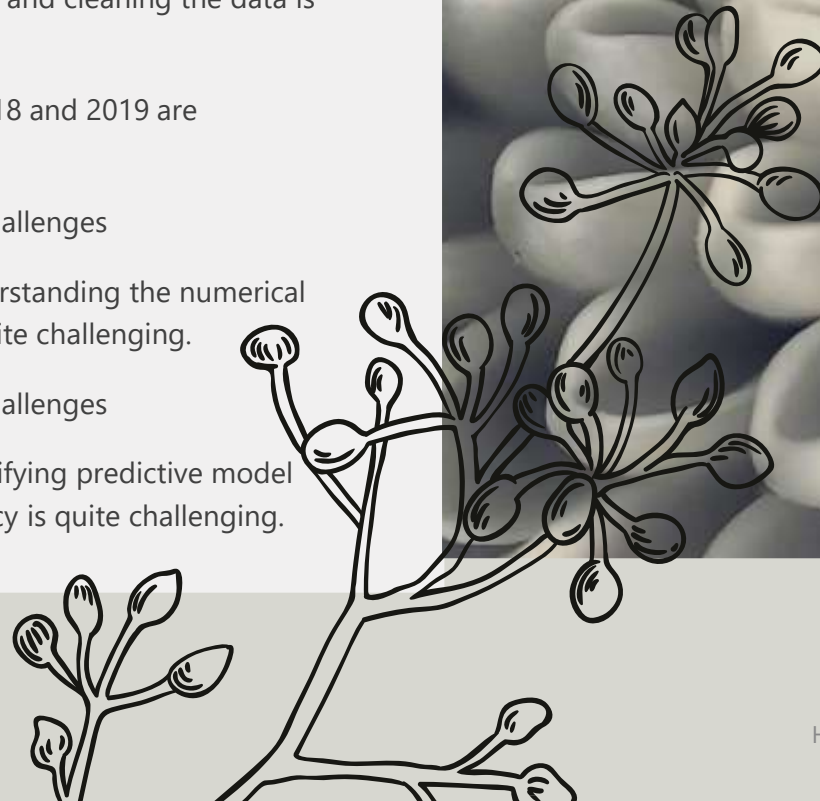
**Update(15'feb):** no challenges

**Update(22'feb):** Identifying predictive model and testing the accuracy is quite challenging.

# Communications

**Weekly status meeting with Professor Rick Lambroff**

**Week – 1 (18'Jan'2023)**

- Professor suggested to use Python for cleaning of dataset instead of doing it manually by Microsoft Excel

- Professor provided tutorial sites for ETL of data processing using Python

**Week – 2 (25'Jan'2023)**

- Professor provided tutorial sites for building a predictive model

- Professor suggested to learn these models and understand clustering algorithms

**Week – 3 (01'Feb'2023)**

- Professor suggested to add more data visualizations after data cleaning process for a better understanding

**Week – 4 (08'Feb'2023)**

- Professor mentioned few changes in the visualizations like adding heatmap, adding same palette colors

**Week – 5 (15'Feb'2023)**

- Professor suggested to try one of the predictive models and test for the accuracy

# Team meetings

| Date | Agenda | Budgeted hours | Attendees | Approval of previous minutes |
|---|---|---|---|---|
| 15/02/2023 | Weekly status update – week 5 | 0.15 | 1. Moganaviniith Rathinavel<br>2. Paras Kishorbhai Gangani<br>3. Ragavi Mudaliyar | Awaiting approval |
| 08/02/2023 | Weekly status update – week 4 | 0.15 | 1. Moganaviniith Rathinavel<br>2. Paras Kishorbhai Gangani<br>3. Ragavi Mudaliyar | Awaiting approval |
| 01/02/2023 | Weekly status update – week 3 | 0.15 | 1. Moganaviniith Rathinavel<br>2. Paras Kishorbhai Gangani<br>3. Ragavi Mudaliyar | Awaiting approval |
| 25/01/2023 | Weekly status update – week 2 | 0.15 | 1. Moganaviniith Rathinavel<br>2. Paras Kishorbhai Gangani<br>3. Ragavi Mudaliyar | Awaiting approval |
| 18/01/2023 | Weekly status update – week 1 | 0.15 | 1. Moganaviniith Rathinavel<br>2. Paras Kishorbhai Gangani<br>3. Ragavi Mudaliyar | Awaiting approval |
| 07/12/2022 | Final group project – submission of SharePoint link, project charter and project proposal | 0.15 | 1. Moganaviniith Rathinavel<br>2. Paras Kishorbhai Gangani<br>3. Ragavi Mudaliyar | Awaiting approval |
| 23/11/2022 | Review of MRP SharePoint Site Follow-up | 0.15 | 1. Moganaviniith Rathinavel<br>2. Paras Kishorbhai Gangani<br>3. Ragavi Mudaliyar | Awaiting approval |
| 16/11/2022 | Review of MRP SharePoint Site Follow-up | 0.15 | 1. Moganaviniith Rathinavel<br>2. Paras Kishorbhai Gangani<br>3. Ragavi Mudaliyar | Awaiting approval |
| 09/11/2022 | Introductory Client meeting - Finalized project topic and dataset | 0.15 | 1. Moganaviniith Rathinavel<br>2. Paras Kishorbhai Gangani<br>3. Ragavi Mudaliyar | Awaiting approval |

# Activities Completed This week

- Collected and securely stored the original data

- Using copies of the original data, clean and prepare the data for    analysis

- The original data is available for the years 2017, 2018 and 2019

- Identifying outliers and data cleaning is completed for the year 2017 using Microsoft Excel

- **Update(25'Jan):** Going through tutorials for ETL of data cleaning instead of manual cleaning is in progress

- **Update(01'Feb):** Completed ETL tutorials and data cleaning for the years 2017, 2018, 2019

- **Update(08'Feb):** Completed data visualization for the year 2019

- **Update(15'Feb):** Completed data visualization for the year 2019, 2018, 2017

- **Update(22'Feb):** Attempted one of the predictive models – K means clustering

# Activities to be Completed Before Next Report

- Preliminary data analysis is to be completed for all the years 2017, 2018 and 2019

- Securely store the cleaned data using naming conventions and version controls

- Identify the databases, languages to be used and develop a functional flow of the project

- **Update(25'Jan):** Data cleaning using ETL python will be completed for all the datasets of years 2017, 2018 and 2019

- **Update(01'Feb):** Understanding predictive models and find a suitable predictive model for our project

- **Update(08'Feb):** Complete the data visualization for all years and start the development of predictive model

- **Update(15'Feb):** Continue development phase II of prediction model

- **Update(22'Feb):** Continue development and testing of predictive models

# Plans for the next phase

**1**

Complete the testing phases successfully

**2**

Include predictive analysis results

**3**

Create interactive dashboard in PowerBI

**4**

Convey the data story to the client through visualizations

# Thank you