

# Infant Mortality Data Analysis; SDG 3: Team Matrix

<b>Anuarg Arora<sup>1</sup></b>	<b>Keshav Gupta<sup>1</sup></b>	<b>Renu Rani<sup>1</sup></b>	<b>Shayan Ray<sup>1</sup></b>
Stony Brook University	Stony Brook University	Stony Brook University	Stony Brook University
SBU-ID: 111425080	SBU-ID: 111464733	SBU-ID: 111482474	SBU-ID: 111424665

## Abstract

In this project, we have analyzed the infant mortality data from Center for Disease Control and Prevention (1, ) and designed a framework to predict the risk of infant death and provide reference to similar pregnancy cases in the past to help the doctors in taking informed decision. Big data technologies such as large scale machine learning, dimensionality reduction, clustering, and similarity search have been employed to make this happen.

## 1 Introduction

Infant mortality is defined as the death of an infant before his or her first birthday. The infant mortality rate is defined as the number of infant deaths for every 1000 live births. The total infant mortality in the year 2016 was 3.15 million. US lags behind other wealthy nations on infant mortality. In 2010, the U.S. infant mortality rate was 6.1 infant deaths per 1,000 live births, and the United States ranked 26<sup>th</sup> in infant mortality among Organization for Economic Co-operation and Development countries. (4, ). Infant mortality rate is an indication of population health, poverty, and socioeconomic status of a country and quality of health services in a country. To ensure development of a country and healthy population, it is important to reduce infant mortality.

In this project, we analyzed the infant mortality data and designed a framework to improve the overall health of a country by predicting the risk of infant death and find similar cases for reference.

Section 2 describes the Sustainable Development Goal and Background of the project. The data-set used in the project is described in section 3. Section 4 goes over the methods used in the project. In Section 5 we share the results of our analysis. Section 6 and 7 discusses the inference and the main contribution of our work. Finally, in Section 8, we enumerate the team member contributions.

## 2 Sustainable Development Goal

The sustainable development goal we aimed to tackle in our project is Sustainable Development Goal 3: *Ensure healthy lives and promote well-being for all at all ages* (3, ). More specifically, we focused on SDG Goal 3.2 which aim to reduce:

1. neonatal mortality to as low as 12 per 1,000 live births and
2. under-5 mortality to as low as 25 per 1,000(bv 40%) live births.

In order to achieve this goal, it would be quite useful to have a framework which can perhaps answer the following questions:

1. Can we predict the risk of infant death?
2. If at risk, can we find similar pregnancy reference cases in the past to aid the medical practitioners in taking well-informed decisions?.

The overall framework proposed uses Spark data pipelines, Dimensionality Reduction, Clustering, Similarity Search, and Large-Scale Machine Learning.

### 3 Data

The data has been taken from Center for Disease Control and Prevention (2, ). Specifically, we are using 2014 "Period Linked Birth-Infant Death Data Files". The data-set contains 4,021,418 observations (split across 3 files). It has approximately 300 features and is over 5 GB in size.

During actual prediction and clustering, the number of features finally used were around 148.

Use the following links to download the raw data (greater than 5GB after zip extraction):

- US-Data
- US-Territories-Data

### 4 Methods

Spark was used as the data pipeline for data loading, data analysis and feature engineering. Following this, the useful features were converted to numeric values, scaled, normalized, transformed with PCA and then used for prediction and clustering. Refer to Figure 1 for an overall summary of methods employed and Figure 2 for detailed spark machine learning pipelines leveraged for this task. Details follow in the sub-sections.

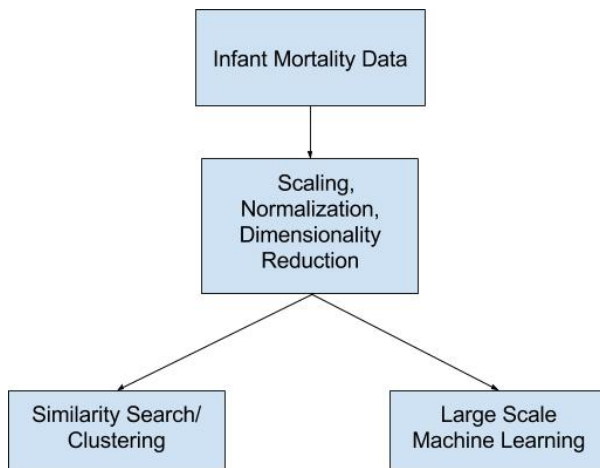


Figure 1: Methods Used

#### 4.1 Spark

Apache Spark is an open source large-scale data processing tool. Spark SQL and dataframes were used for data processing. Specifically, pyspark.ml library was used for feature engineering, dimensionality re-

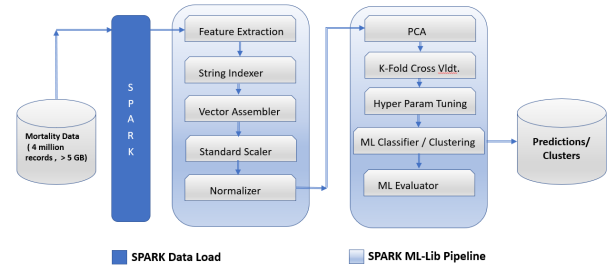


Figure 2: Implementation Framework

duction, clustering and for various machine learning models.

#### 4.2 Feature Engineering

Feature encoding, scaling and normalization of features are performed in order to standardize the data and reduce the frequency of out-liers. In addition to this, PCA (Principal Component Analysis) is used to perform dimensionality reduction for reducing the high-dimensionality/ number of features of our data. The advantage of reducing the number of features (from 148 to 50 in our case) would be avoiding over-fitting. It leads to better prediction accuracy and a more generalized prediction model.

#### 4.3 Exploratory Data Analysis

Some insights gained from data distribution were as follows:

1. Correlation among some key features: (Refer to Figure 2). It became evident from data that the infant death risk was not a clear correlation of a few factors but a complex mix of many parameters.
2. Fetal cause of death distribution: (Refer to Figure 3) The leading cause of death based on the ICD10 (International Cause of Death) was extreme pre-maturity(code=P072).
3. Distribution of some key features in the dataset: (Refer to Figure 4).

#### 4.4 Machine Learning

We used supervised learning algorithms for doing our analysis. The label or the ground truth was the

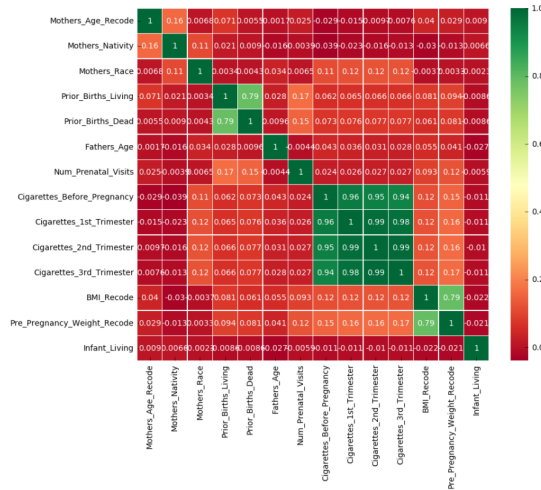


Figure 3: Correlation among few key features

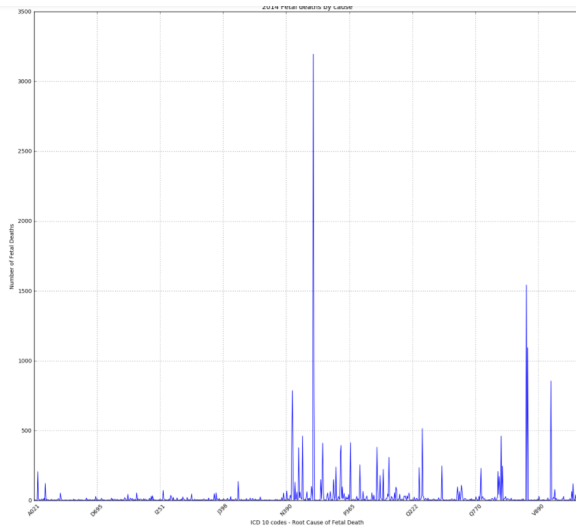


Figure 4: Root Cause of Fetal Death

information whether the infant was living as indicated by a row of data. Other features such as those indicating the Five Minute APGAR Score, health of mother were taken as the features for the analysis. We used the following machine learning methodologies:

1. Logistic Regression
2. Gradient Boosted Trees
3. Random Forest
4. Decision Tree Classifier

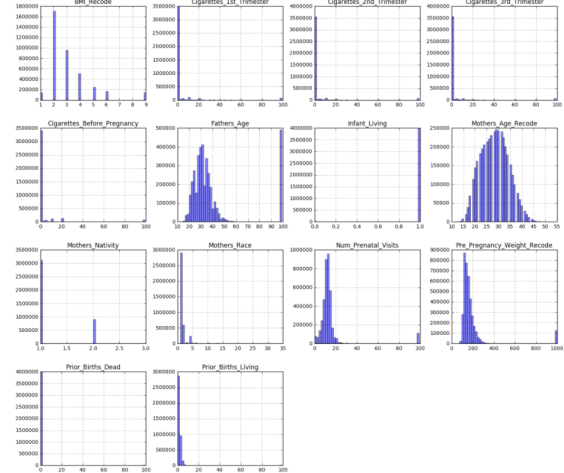


Figure 5: Sample Feature Distributions

## 4.5 Similarity Search

We used similarity search to find the similar cases which can aid the medical practitioner to take well informed decision. Cosine similarity is calculated on dimensionality reduced data i.e after applying the PCA on the observations. Also, before finding the similarity, we are taking a random sample from all the observations and then we apply the similarity search on the randomly chosen sample.

## 4.6 Clustering Algorithm

K-Means [3] clusters the observations in multiple clusters (k=20 in this case) such that the observations belong to the cluster with the nearest mean. For any new observation that arrives, it is clustered to the most similar existing cluster. It enables us to find similar reference cases in the past.

## 5 Results

All the results in HTML format are submitted along with this report archive. Refer to it for detailed dataset analysis and respective results. The following Machine Learning algorithms were employed:

1. Logistic Regression
2. Gradient Boosted Trees
3. Random Forest
4. Decision Tree Classifier

Logistic Regression	
AUROC	0.5
AUPRC	0.50289
Gradient Boosted Trees	
AUROC	0.99999
AUPRC	1.0
Random Forest	
AUROC	0.99999
AUPRC	0.99989
Decision Tree Classifier	
AUROC	1.0
AUPRC	1.0

**Table 1:** Machine Learning Results

The results for the Machine Learning algorithms are presented in Table 1. AUROC (Area under the Receiver Operating Characteristic) and AUPRC (Area under the Precision Recall Curve) were used to evaluate the performance of these models. Refer to for more details on AUROC and AUPRC.

In general, the closer these values are to 1, the better the accuracy of the model. From the results it was observed that a naive model like Linear Regression started off with values of 0.5 (AUROC and AUPRC) whereas more powerful models like Decision Tree Classifier and ensemble models like Random Forests performed way better with very high accuracy (1).

TODO: Add about similarity search and explain the table.

## 6 Discussion

For a pregnant mother, it is possible to predict the risk of infant death with this framework and US CDC data. If a risk is detected, the medical practitioners can leverage this framework to find the similar pregnancy cases in the past which are a close match to this one and perhaps refer the pregnant mother to a Perinatologist with proven track record (specialist in mother and fetal health) instead of a regular gynecologist. This would definitely reduce the chances of infant mortality, thereby fulfilling SDG 3.2 to an extent. After finding that there is a high risk to the pregnant mother and to her fetus, similarity search and clustering can be used to find similar cases. These similar cases will help the medical practitioner to make a well informed decision by looking at what type of treatment might have

worked in past or what type of issues can occur in future for the current case.

## 7 Conclusion

This framework serves as a tool to answer the following question:

1. Can we predict the risk of infant death? Yes with the help of large-scale machine learning models, the risk or no-risk classification has been successfully accomplished.
2. If at risk, can we find similar pregnancy reference cases in the past to aid the medical practitioners in taking well-informed decisions? Yes, using clustering for big-data and similarity search for sampled data, this has been successfully accomplished.

This framework empowers the medical practitioners to assess the risk of infant death ahead of time and look-up past matching pregnancies as treatment reference. Perhaps this can lead to the pregnant mother's case being transferred to a Perinatologist instead of a regular Gynecologist

Additionally, TODO: ADD ABOUT SIMILARITY SEARCH: TODO

## 8 Team Member Contributions

Data Load: Shayan, Anurag  
Data Cleaning: Renu, Keshav  
Feature Extraction: Keshav, Renu  
Dimensionality Reduction: Keshav, Anurag  
Clustering: Shayan, Anurag  
Machine Learning: Shayan, Renu  
Similarity Search: Anurag, Keshav  
Model Evaluation: Shayan, Renu

## References

- [1] CDC website
- [2] CDC Data
- [3] UN Sustainable Development website
- [4] K-Means Clustering Wikipedia
- [4] National Vital Statistics Reports (NVSS) Report on Infant Mortality