# Understanding secondary factors contributing to child mortality

CSE 545 - Big Data Analytics (Spring `20)
Team XYZ - Karan Shah, Komal Agarwal, Narayan Acharya, Vikas Dhyani
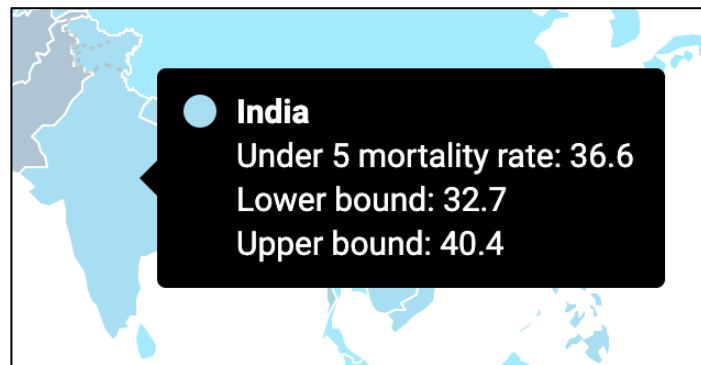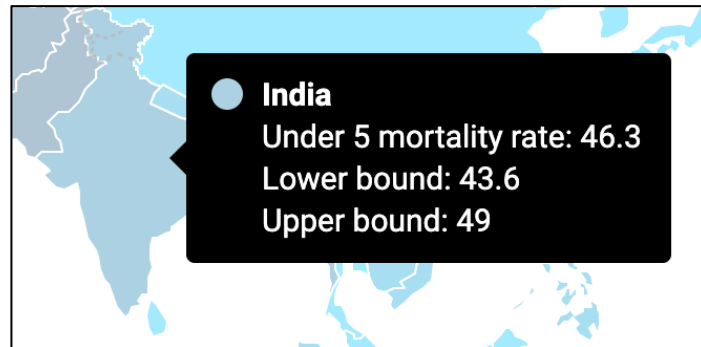
# What are we doing?

We look to **study the impact of secondary factors** that can help decrease child mortality rate: Environment, Amenities, Lifestyle, Health. And**,** we try to **suggest priority actions** for Government bodies to action.

# Why should one care?

- "…**52 million children under 5 years of age, will die between 2019 and 2030**." *(Source: UNICEF)*
- Even though we are seeing progress it is NOT enough to control the disturbing projections we see for the decade.
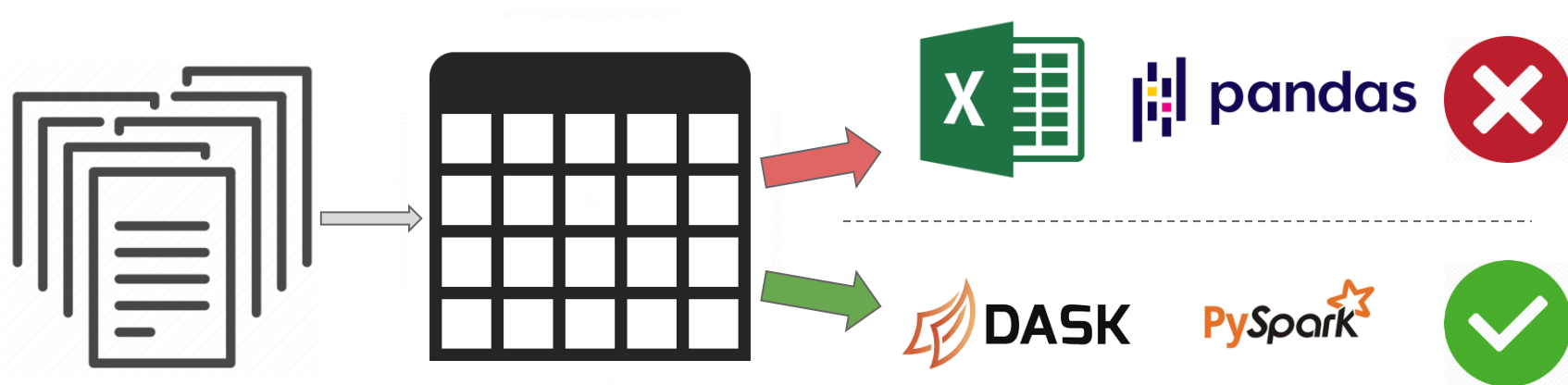
# How does it relate to SDG?

- SDG 3 -- **"Ensure healthy lives and promoting well-being for all at all ages."**
- Containing infant mortality helps mental and physical well-being of mother and child helping them live healthy lives.



**India**
Under 5 mortality rate: 46.3
Lower bound: 43.6
Upper bound: 49

**India**
Under 5 mortality rate: 36.6
Lower bound: 32.7
Upper bound: 40.4

*Under 5 mortality rates for India in 2014 (top) vs 2018 (bottom) from UNICEF data.*
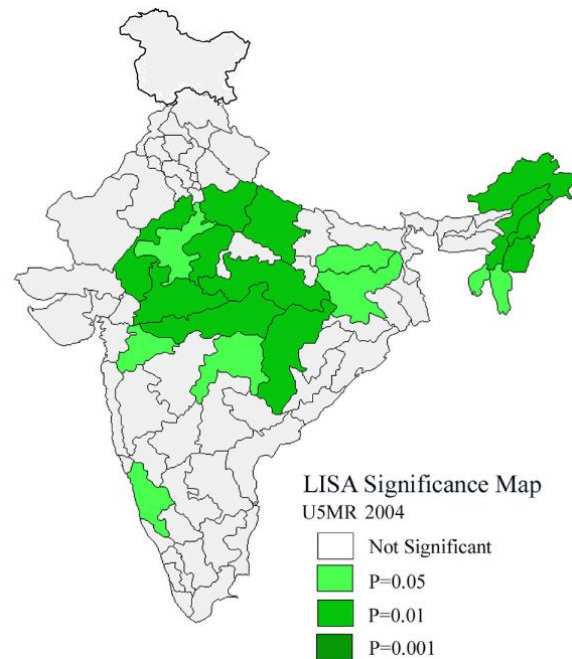
# Why Big Data?



RAW DATASET:
16.81 GB

> 3M Records
> 200 Columns/Record

Traditional Programs suffer when entire data is not in memory.

Traditional Programs do NOT support distributed and parallel processing out of the box.

# Background

- A Geospatial Analysis [study](2011) has shown intra-state and inter-regional disparities in infant mortality in India with help of [geospatial techniques](#) like,
  - Moran's-*I*
  - univariate LISA
  - bivariate LISA
  - Spatial error regression
  - spatiotemporal regression



LISA Significance Map
U5MR 2004

☐ Not Significant
■ P=0.05
■ P=0.01
■ P=0.001

LISA ([Cluster and Significance](#)) map depicting spatial clustering and spatial outliers of under-five mortality by incidence of poverty across 74 geographic regions

# Background

**What is exactly a Infant Mortality Rate?**

Number of deaths of infants (age under one year) in 1,000 <u>live births</u> in a year.

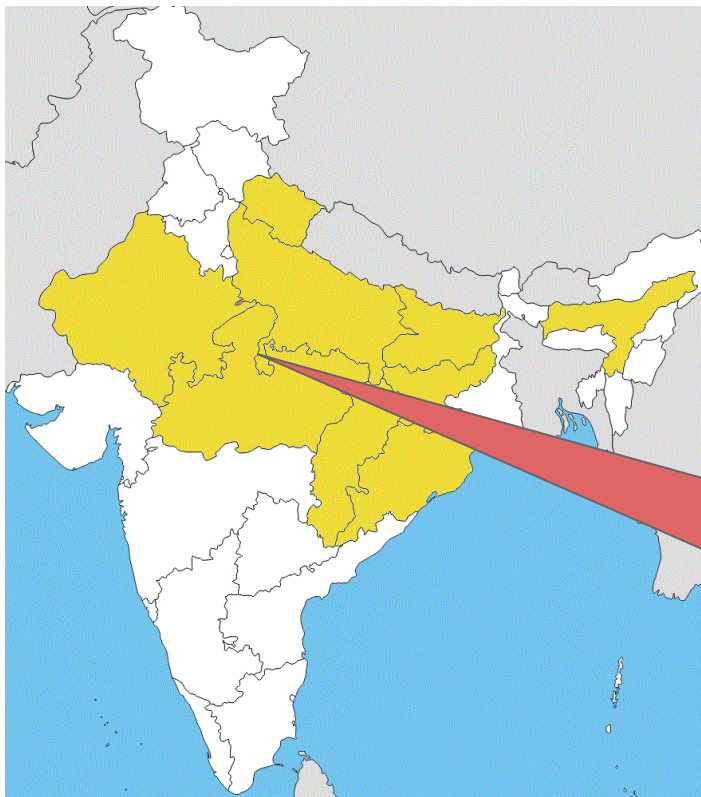**Live births? How is it different than just…. "births"?**

Stillbirth - death or loss of a baby before or during delivery

Live birth - baby's birth, showing any sign of life after exiting maternal body

# Data



children_lost

(label)

=

born_alive_total - surviving_total

(features)

- 9 states of India
- ~16.81 GB size
- 200+ features (environment, amenities, lifestyle, health,..)
- 3M records (each record for a woman, district wise)

A Digital India Initiative

data.gov in
Open Government Data (OGD) Platform India

🏠 / Annual Health Survey : Woman Schedule

## Catalog Info

Get data on Annual Health Survey : Woman Schedule. Woman Schedule comprised two sections. Section-I contains information relating to the outcome of pregnancy(s) (live birth/still birth/abortion); birth history; type of ...   [+]

O Released Under: National Data Sharing and Accessibility Policy (NDSAP)

O Contributor:
  Ministry of Health and Family Welfare
  Department of Health and Family Welfare

O Keywords  AHS   woman   Pregnancy
  Abortion   Delivery   Natal
  Immunization   Children   Supplement
  child   Pneumonia   Diarrhoea
  < Show more >

O Group: Annual Health Survey

O Sectors:  Family Welfare   Health
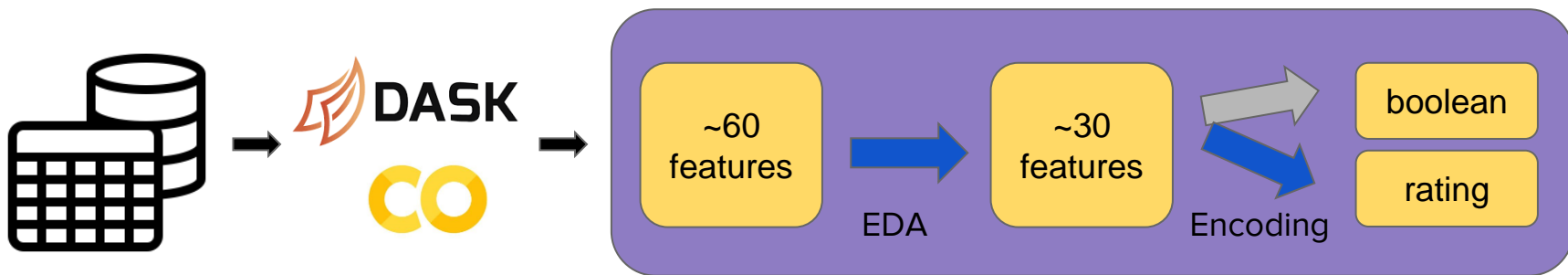  Health and Family welfare

O Published on Data Portal: 20/08/2019

O Source: Open Government Data (OGD) Platform India

3 GOOD HEALTH AND WELL-BEING

# Method

1) **Preprocessing, EDA**

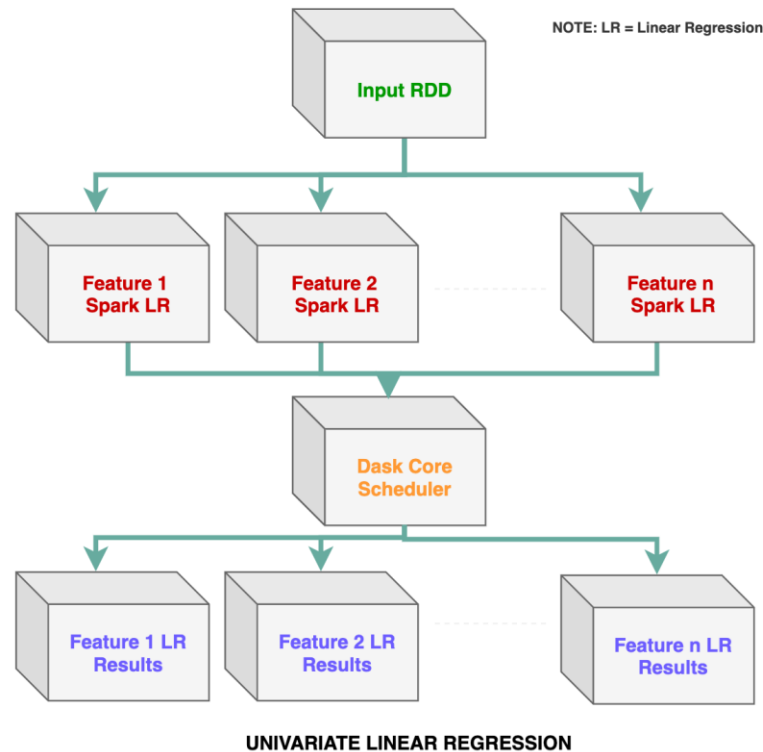   ○ Based on EDA on 1 state data using Google Colaboratory,



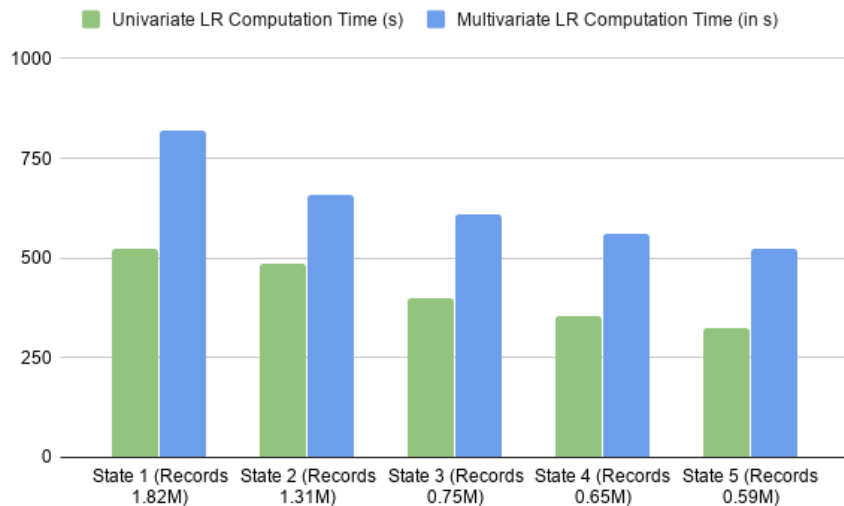   ○ Pipeline created to generate full data on GCP cluster from all 9 states raw data

# Linear Regression

- We had close to 30 features for independent linear regressions.
  - Dask for parallel processing of multiple independent univariate linear regressions.
  - Individual linear regressions were computed using Spark.

- We used only Spark for Multivariate Linear Regression.
  - Since there was only one linear regression computation, Dask was not needed.



NOTE: LR = Linear Regression

Input RDD

Feature 1 Spark LR | Feature 2 Spark LR | Feature n Spark LR

Dask Core Scheduler

Feature 1 LR Results | Feature 2 LR Results | Feature n LR Results
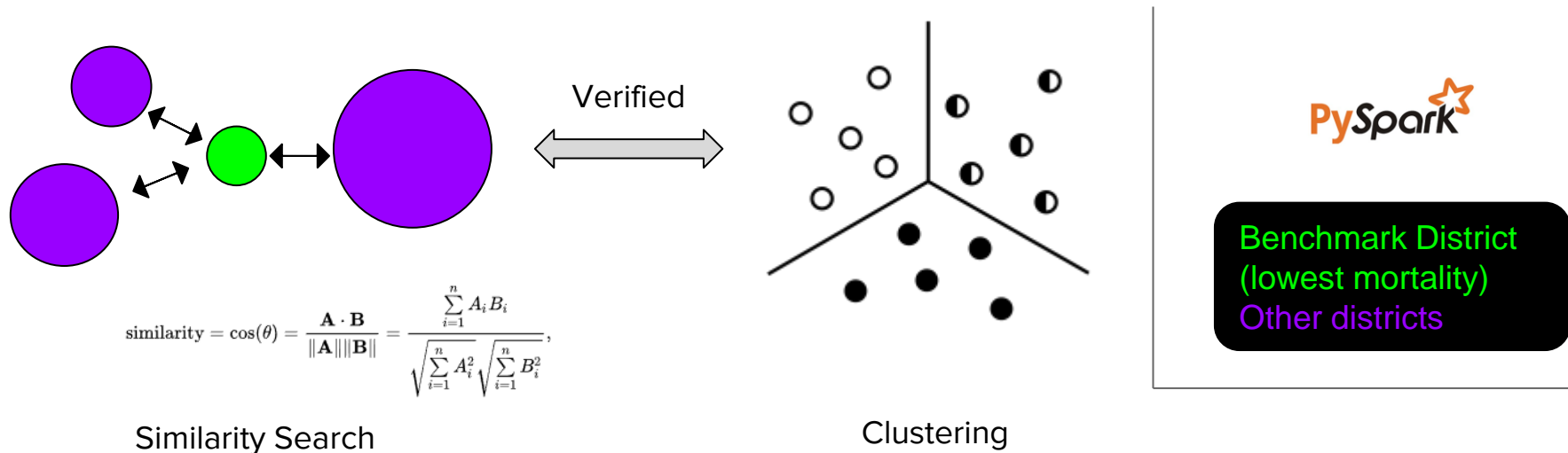
UNIVARIATE LINEAR REGRESSION

# Hypothesis Testing & Computation Time

- We select the top 5 positively correlated features to then calculate p-values.
- We correct for multiple features in Multivariate Linear Regression using Bonferroni Correction.

- Due to the number of records (more than 1M for a state), we use Stochastic Gradient Descent with a learning rate of 0.1 to learn the beta coefficients.
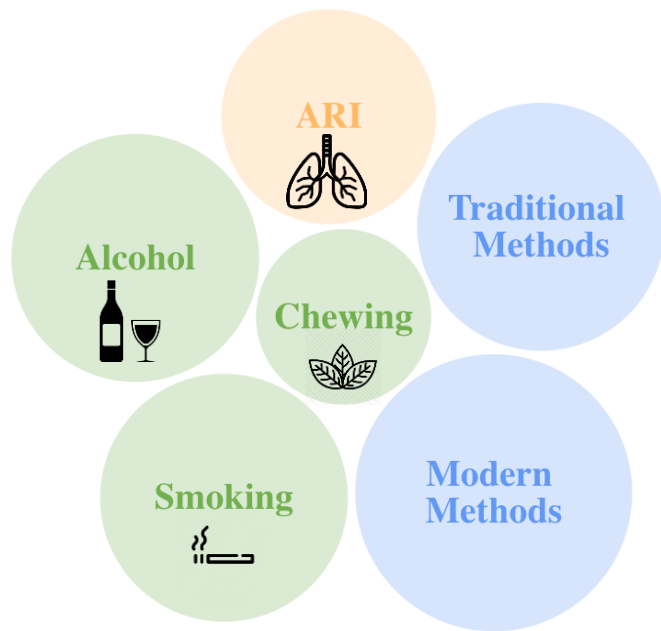


*Computation Time for 5 states*

# Similarity Search & Clustering



Verified

$$similarity = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum\limits_{i=1}^{n} A_i B_i}{\sqrt{\sum\limits_{i=1}^{n} A_i^2} \sqrt{\sum\limits_{i=1}^{n} B_i^2}},$$

Similarity Search

Clustering

PySpark

Benchmark District
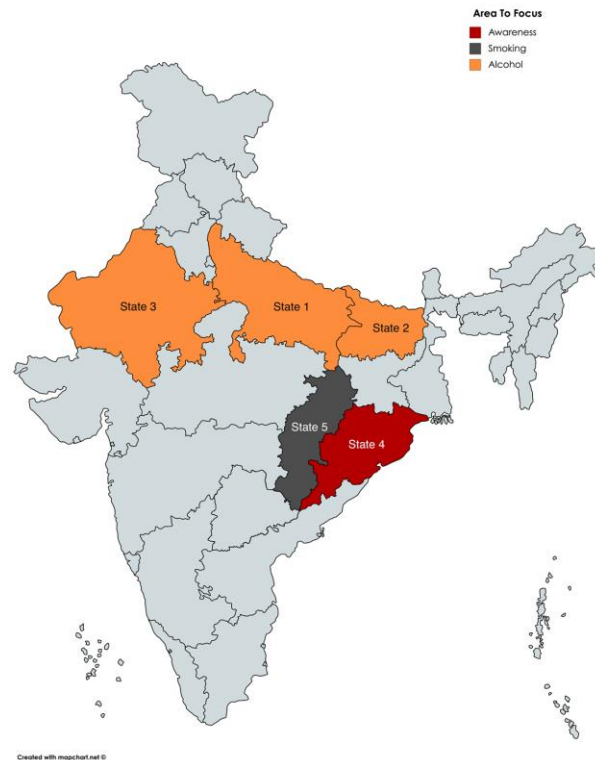(lowest mortality)
Other districts

- Similarity of each district was calculated with benchmark district.
- It was identified that districts with low literacy, no amenities indeed have high mortality.

Clustering based on chosen features verified that the clusters included similar districts; we chose number of clusters to be 4.
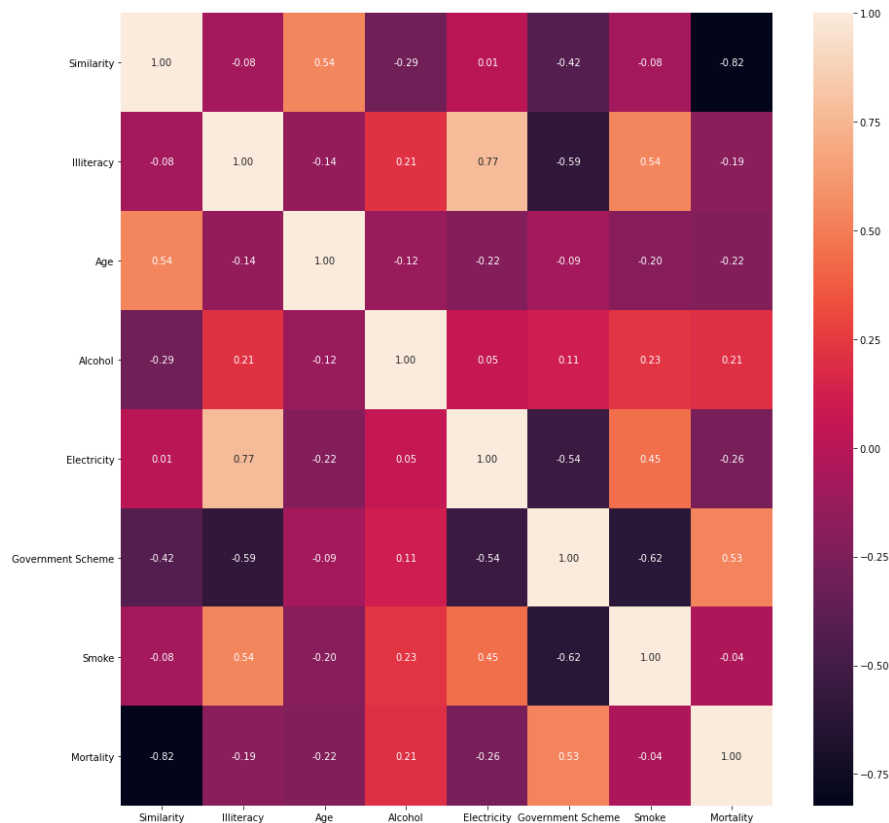
# Linear Regression Results



Top Features That Correlated Most With Mortality



Top Feature To Work On For State
(Ignoring Modern & Traditional Methods)

# Similarity Search Results



Correlation between Non-biological factors and Mortality…
- Similarity (of districts with benchmark)
- illiteracy
- Age
- Alcohol
- Electricity
- Government Scheme
- Smoke
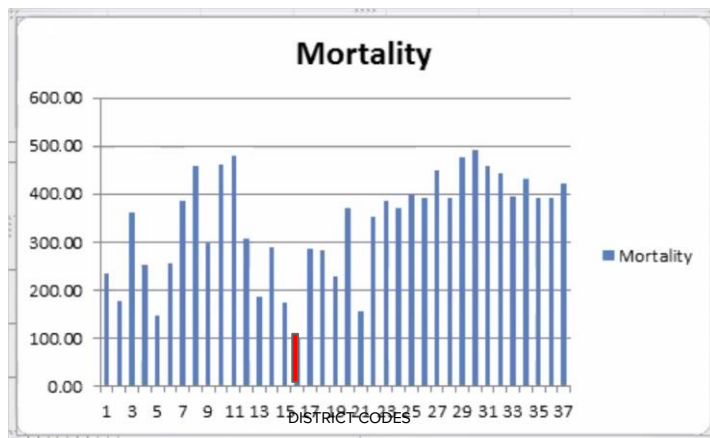
# Similarity Search Results (2)



Fig2. Mortality Rate for every district in bihar

*Mortality Rate For State 2 Bihar*
*For Every District*
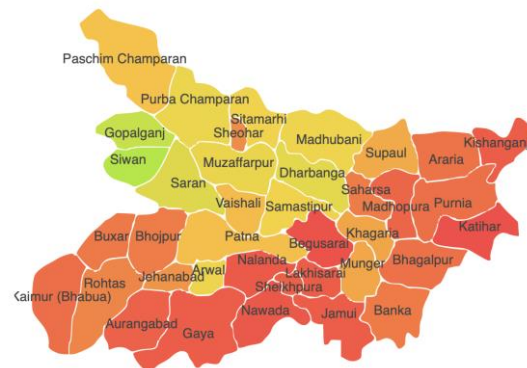*Highlighted District is our Benchmark District - Siwan*



Fig1. Similarity for every district in bihar

*Similarity For Every District with Benchmark District*
*Siwan For State 2 (Bihar)*

# Conclusion

- Considering SDG 3 -- "Ensure healthy lives and promoting well-being for all at all ages.", our results conclude that non biological factors <u>also</u> impacts IMR at some extent.

- The analysis reveals actions and awareness measures towards non biological factors such as Education, Alcohol consumption, Smoking habits etc. should be taken to reduce IMR.