

A STORY THROUGH DATA

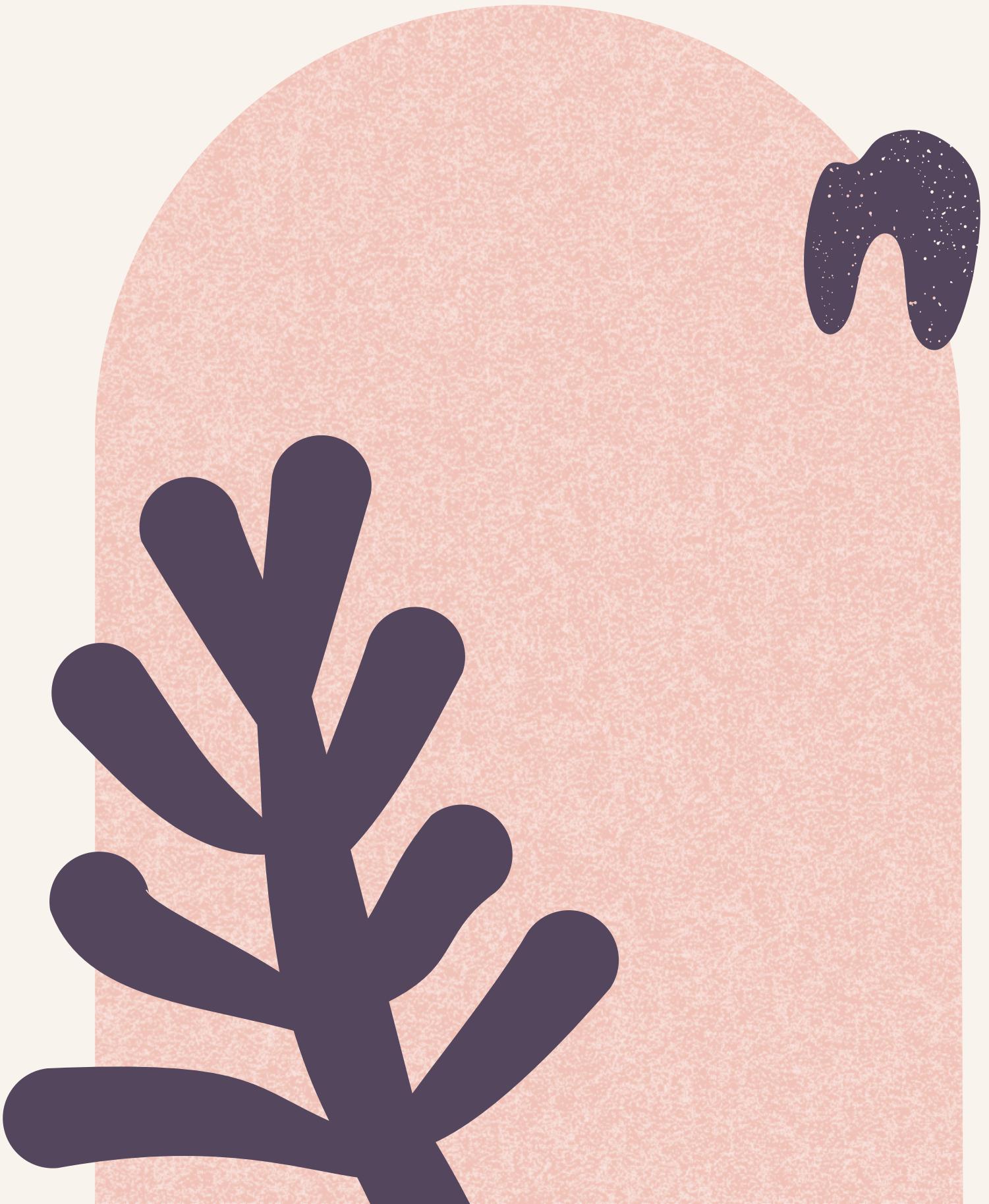
Predicting Health Disparity in the United States

ELINA RANKOVA

DATA SCIENTIST
BROOKLYN, NY

Presentation Agenda

- BUSINESS OBJECTIVE
- DATA UNDERSTANDING
- MODELING
- RESULTS
- LOOKING TO THE FUTURE
- CONTACT INFO



Business Objective

In recent years more companies and institutions are specializing in public health are leaning on data and technology to inform business decisions.

This project aims to identify overall health disparity across the United States by defining a Health Disparity Index to help direct business opportunities and close much needed care gaps.

Data Understanding

About the Data

CDC PLACES and SDOH data spanning 2017-2021
8287 rows from the SDOH data and 780890 from the PLACES

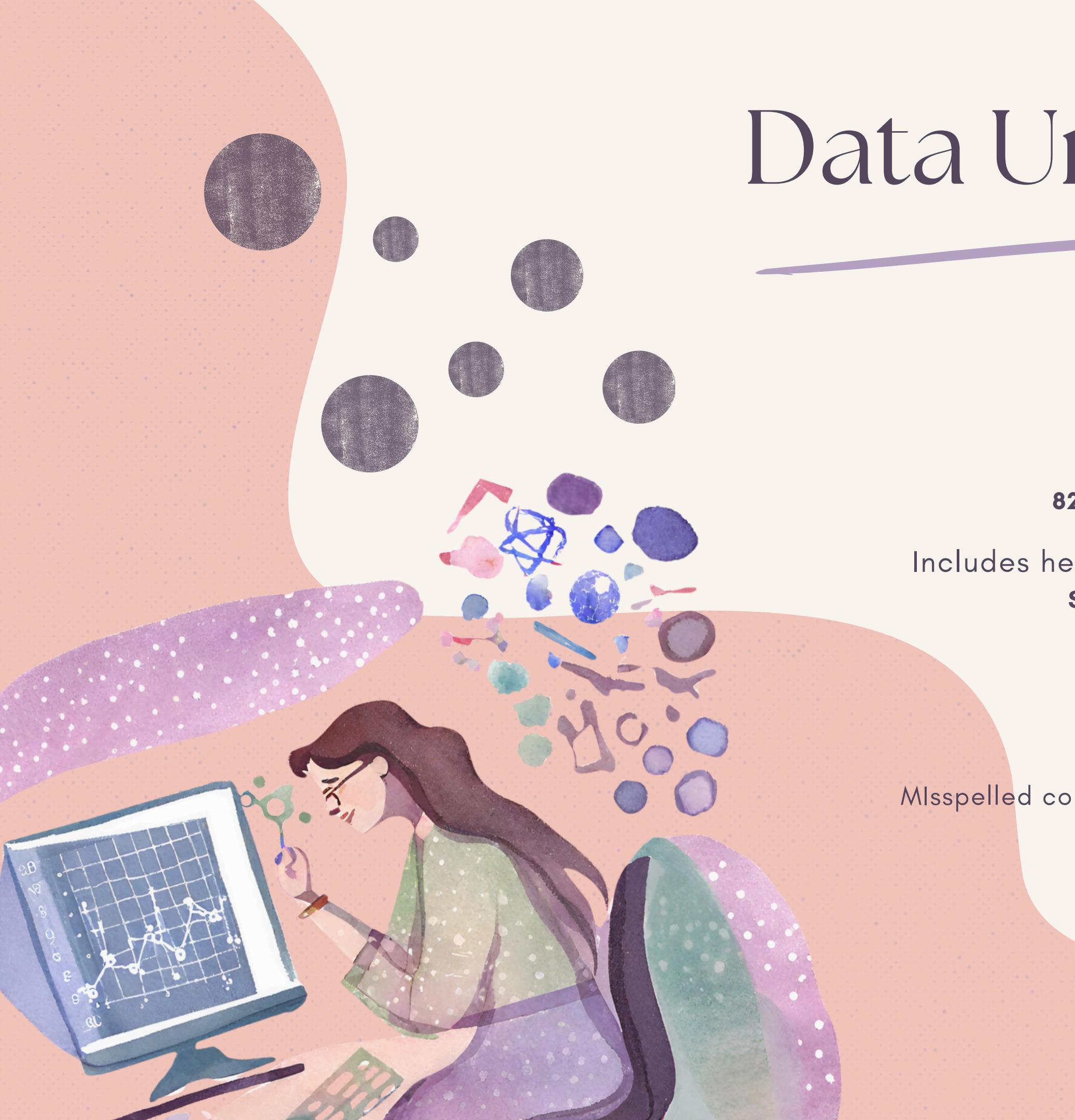
Includes health measure and categories for counties in the US
SDOH, Health Outcomes, Prevention, Health Risk Behaviors, Health Status, Disabilities

Columns with too much missingness are dropped
Columns not present in SDOH were dropped as well

Misspelled columns were combined with the misspelled column dropped
`Geolocation`/`Geolocation` & `MeasureID`/`MeasureId`

Proper pre-split transformations applied
`State` → categorical, `LocationID` → object

Features with either duplicated or unnecessary information were dropped



Data Understanding

Category
measure group

Feature Preview

State
two letter
abbreviation

Data Value Type
percentage aggregation type

Location ID
county unique identifier

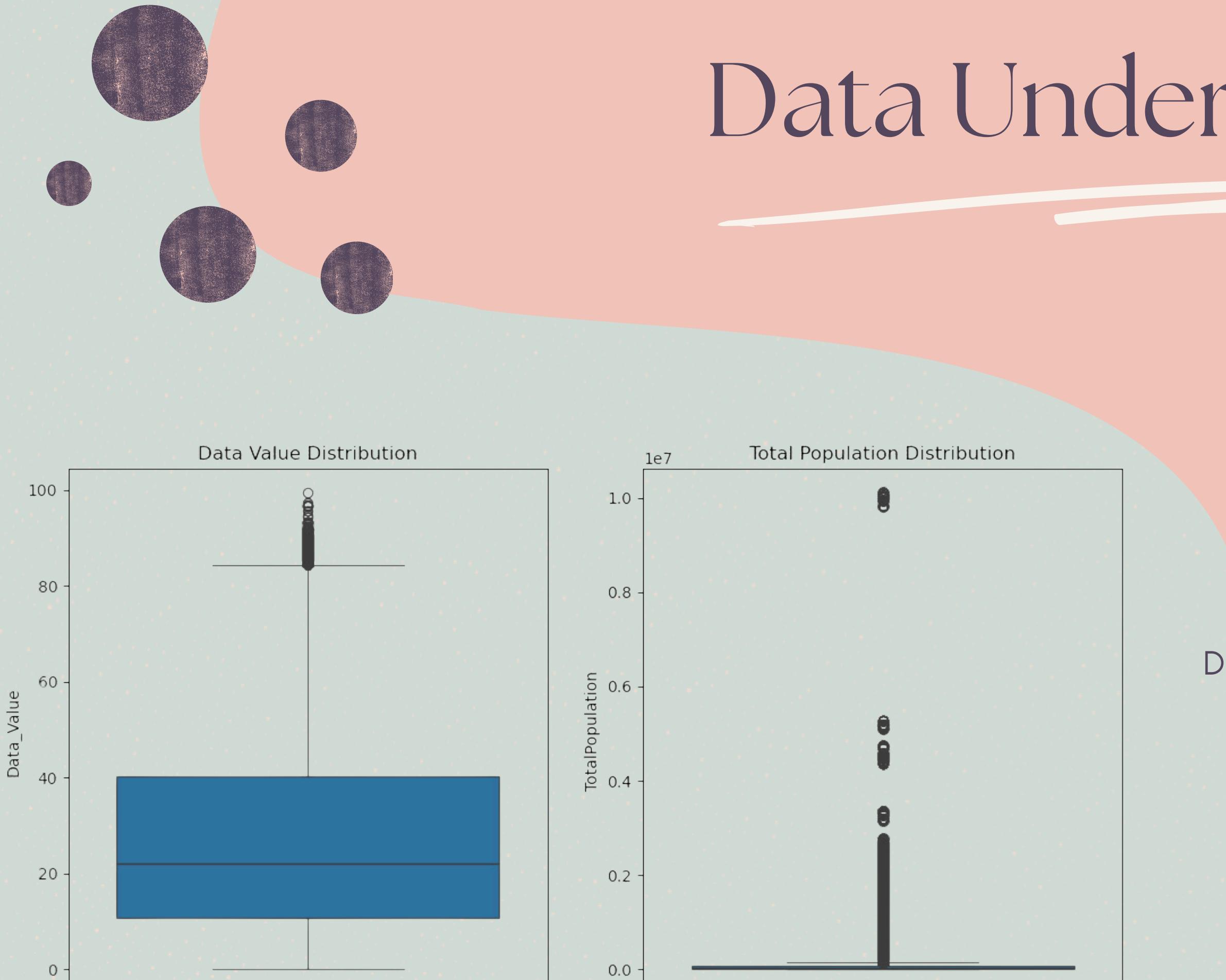
Location Name
county name

TotalPopulation
population in the county

Short Answer Text
short text of measure

Data Value
percentage of measure
depending on data value
type

Data Understanding



Statistical Analysis

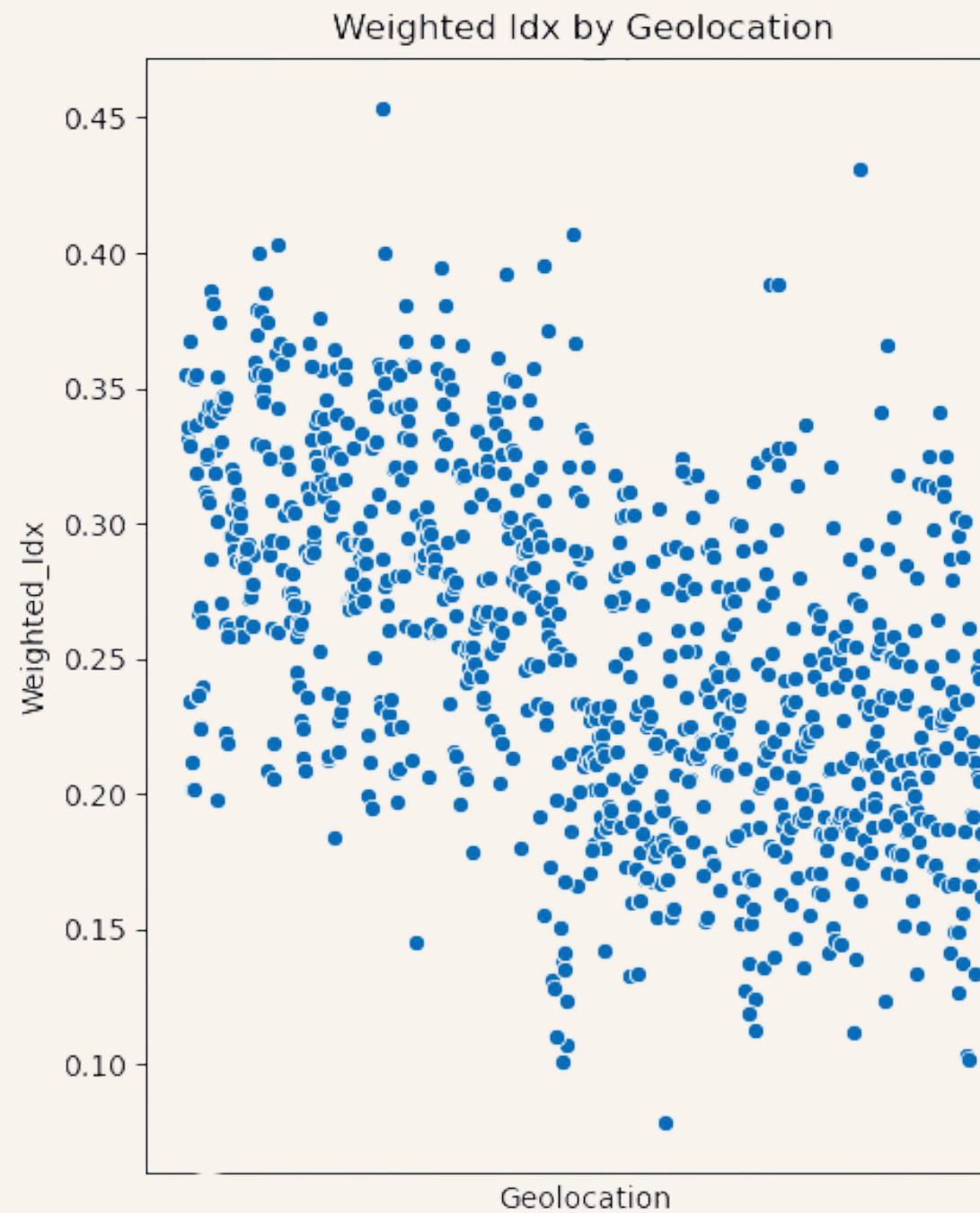
Outliers were dropped using IQR
Data_Value and TotalPopulation

Data_Value was normalized with new
Scaled_Value feature created
**`RobustScaler` used to keep
consistent with IQR analysis**

Data Understanding

Health Disparity Index

Scaled data values aggregated by Geolocation and Data Value Type and weighted by the population.

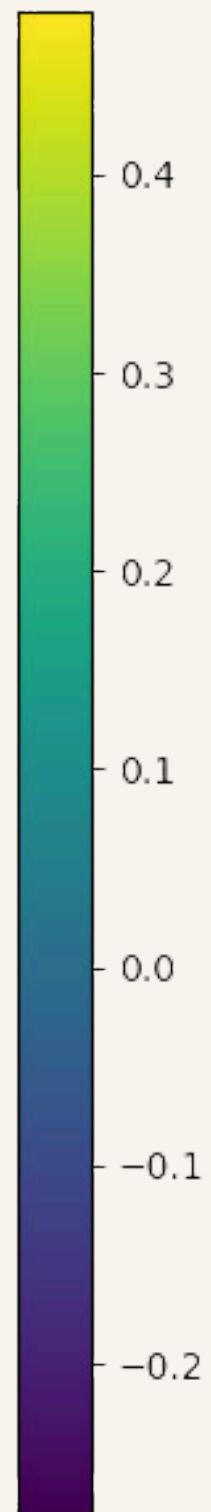
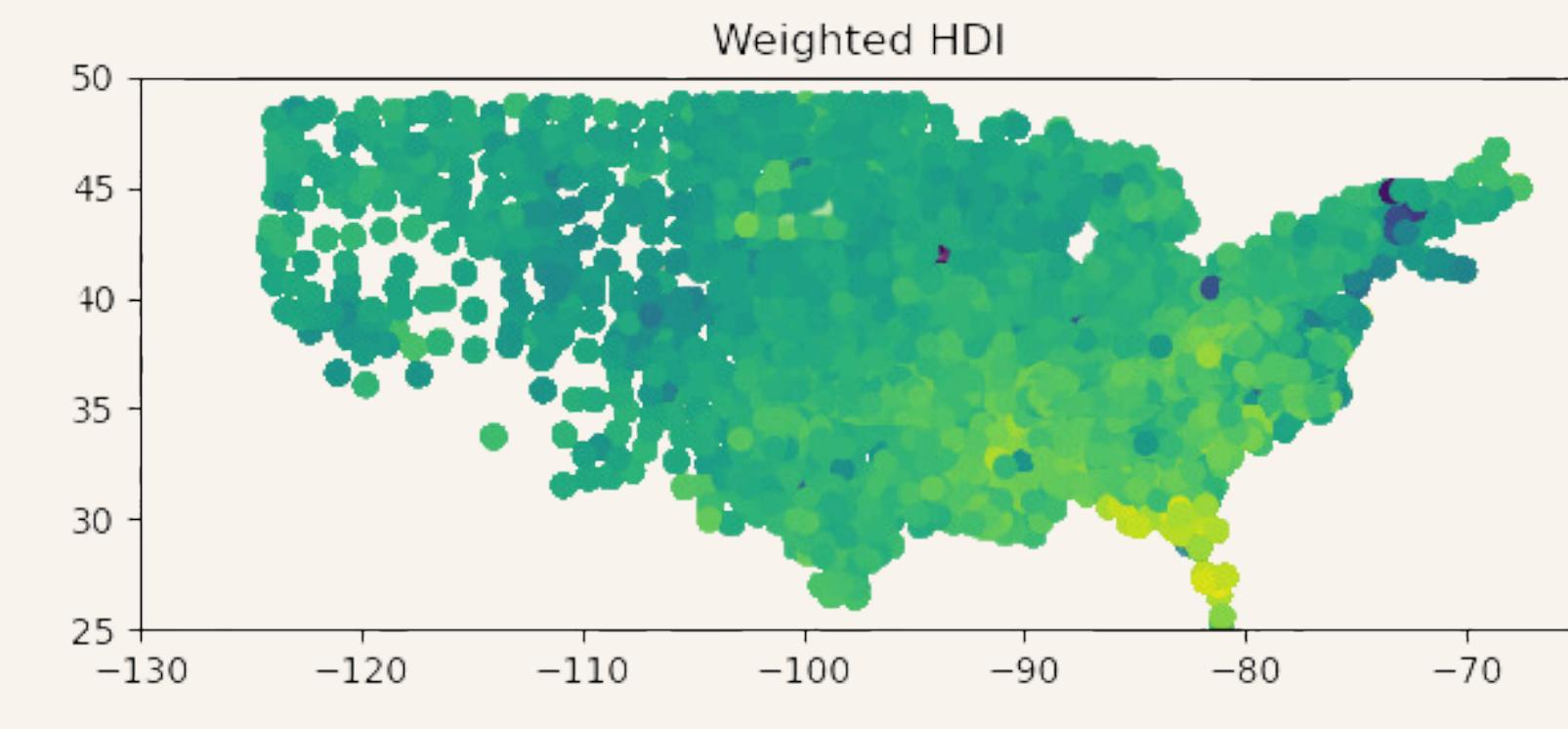
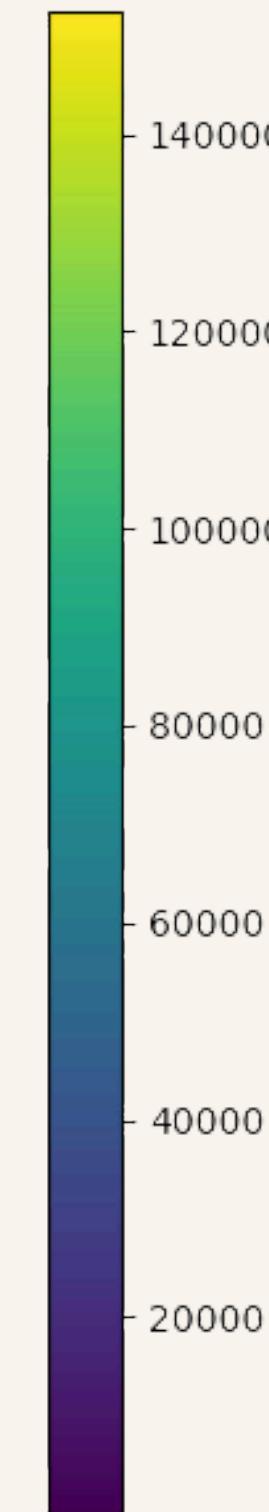
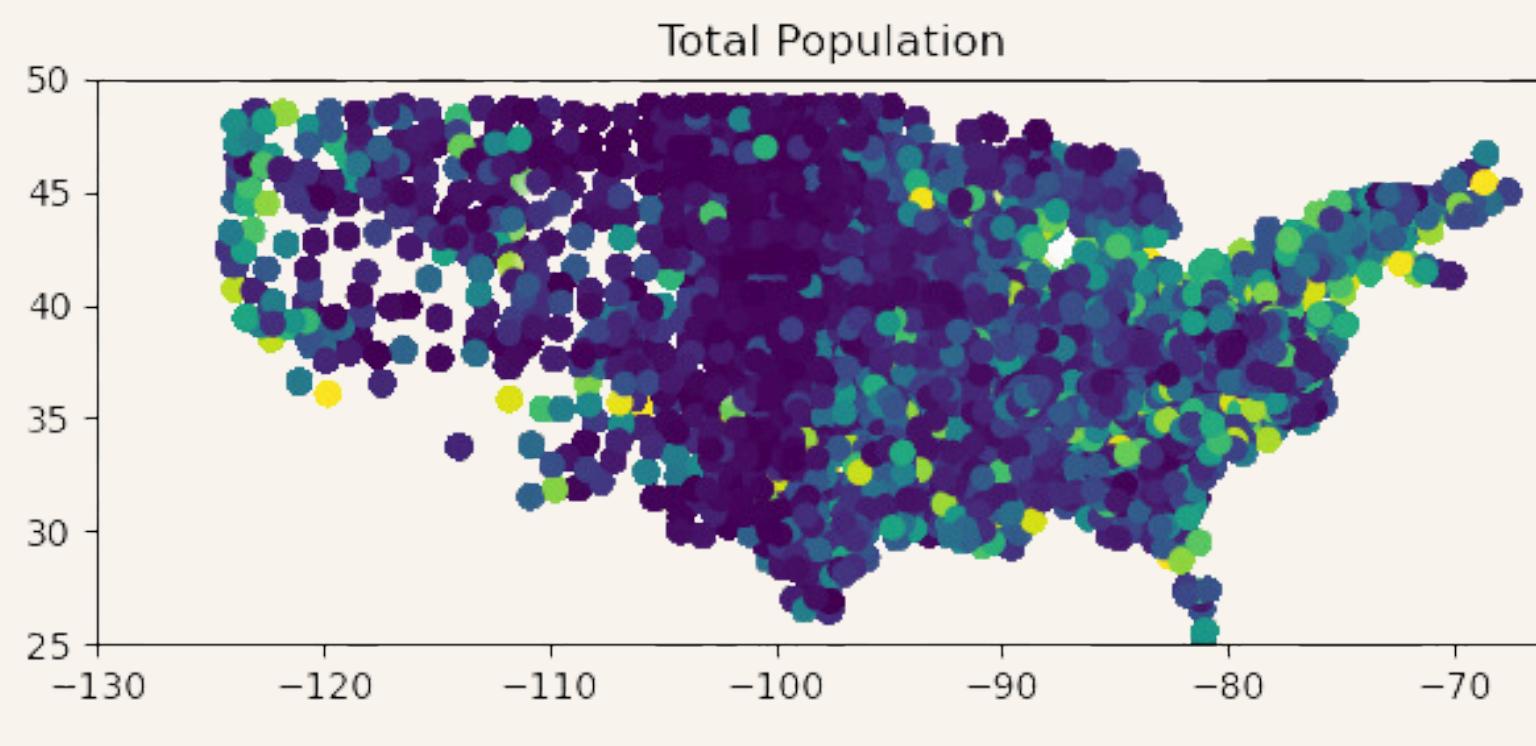


All aggregation was done on a record level so that each scaled value was aggregated based on the record's health measure and data value type.

When population weight is added, we see geographical distinction

Data Understanding

Health Disparity Index



Modeling Base



METRIC	TRAIN VALUE	TEST VALUE
RMSE	0.036	1.68E+09
R-SQUARED	0.697	-6.707E+18
MAE	0.032	1.909E+06



Cross validation and parameter tuning is conducted through a custom class offering the option of `RandomizedSearchCV` or `gp_minimize` for parameter tuning

Modeling

Lasso

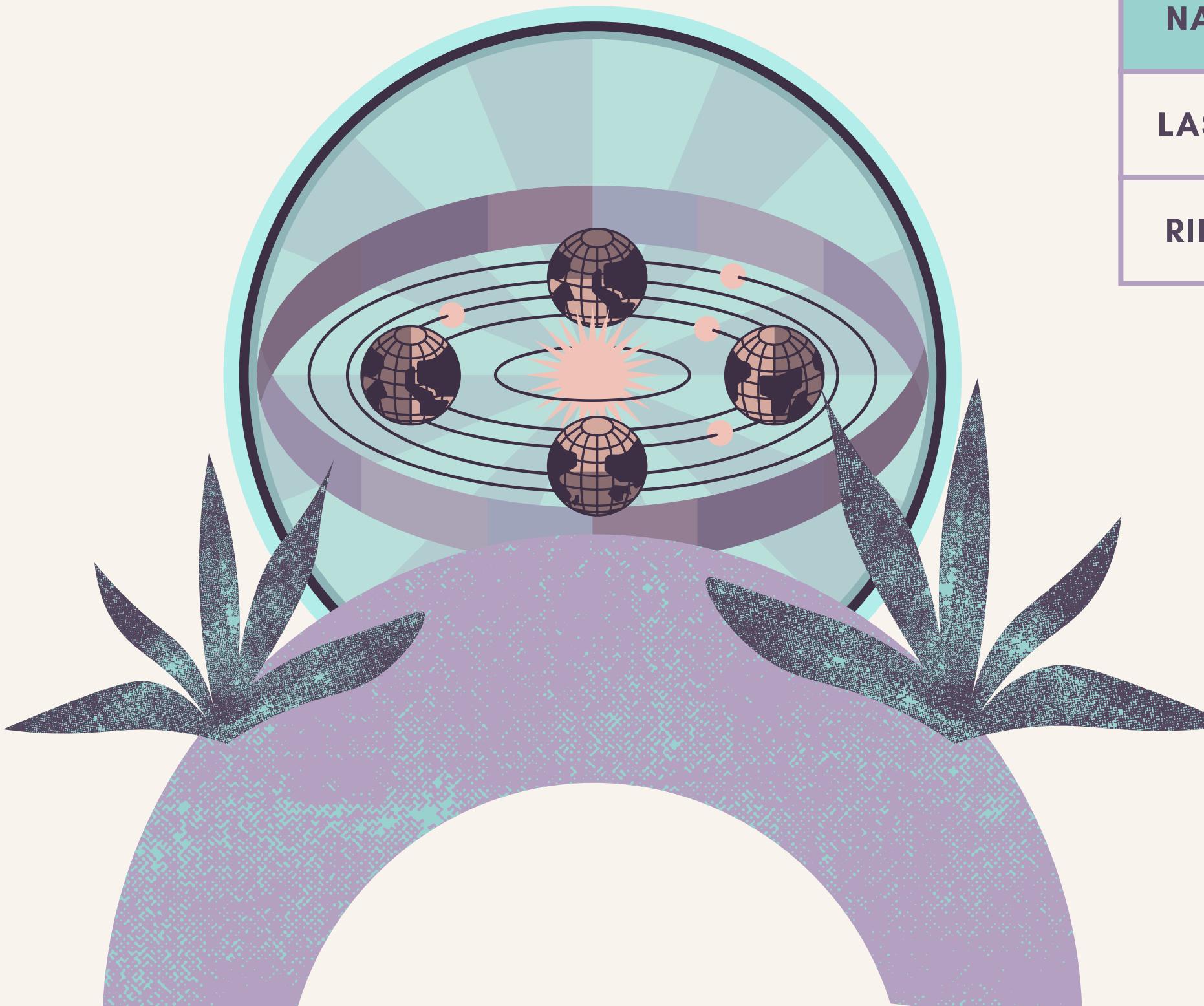


MODEL NAME	TRAIN SCORE	TEST SCORE	TRAIN R2	TEST R2	TRAIN MAE	TEST MAE
LASSO	0.059	0.059	0.181	0.180	0.048	0.048



Modeling

Ridge



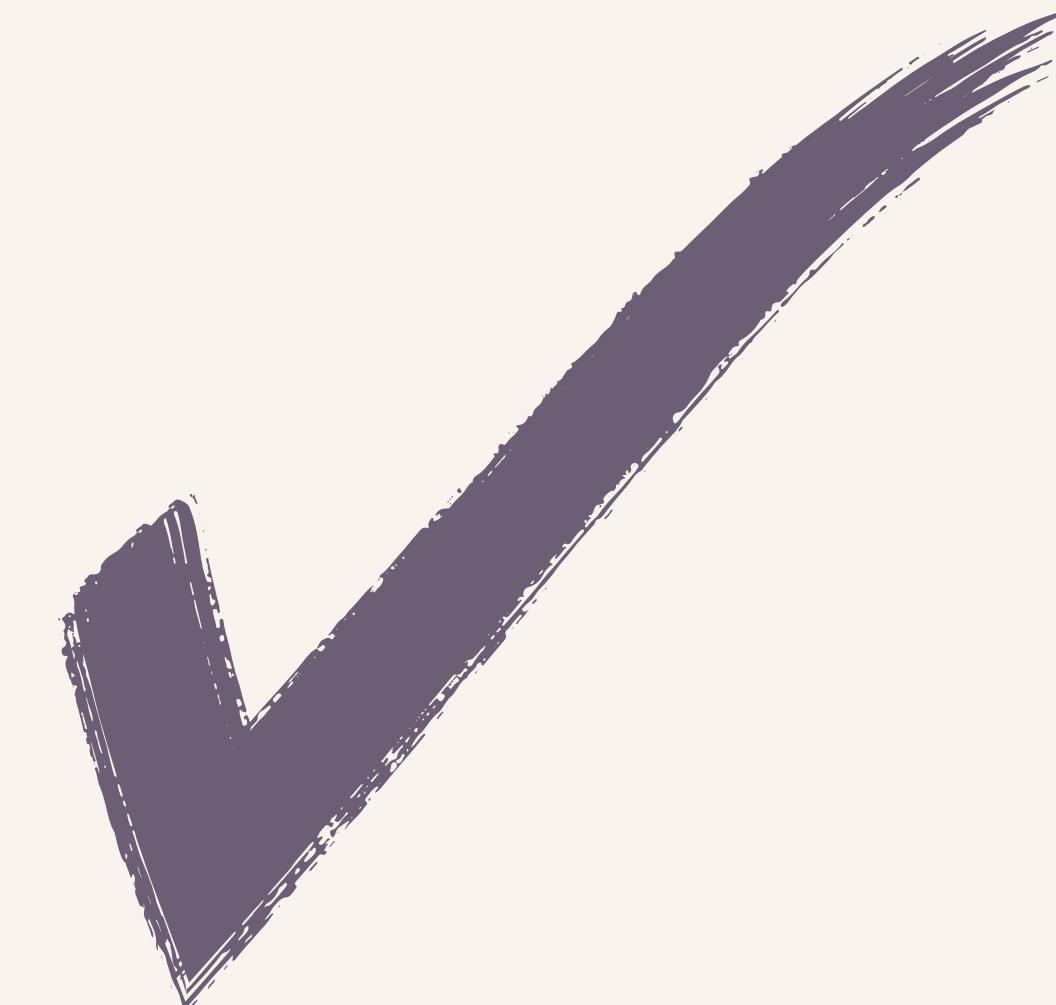
MODEL NAME	TRAIN SCORE	TEST SCORE	TRAIN R2	TEST R2	TRAIN MAE	TEST MAE
LASSO	0.059	0.059	0.181	0.180	0.048	0.048
RIDGE	0.036	0.037	0.697	0.680	0.032	0.032

Modeling

Ridge



MODEL NAME	TRAIN SCORE	TEST SCORE	TRAIN R2	TEST R2	TRAIN MAE	TEST MAE
LASSO	0.058	0.059	0.181	0.180	0.043	0.048
RIDGE	0.036	0.037	0.697	0.680	0.032	0.032



Modeling

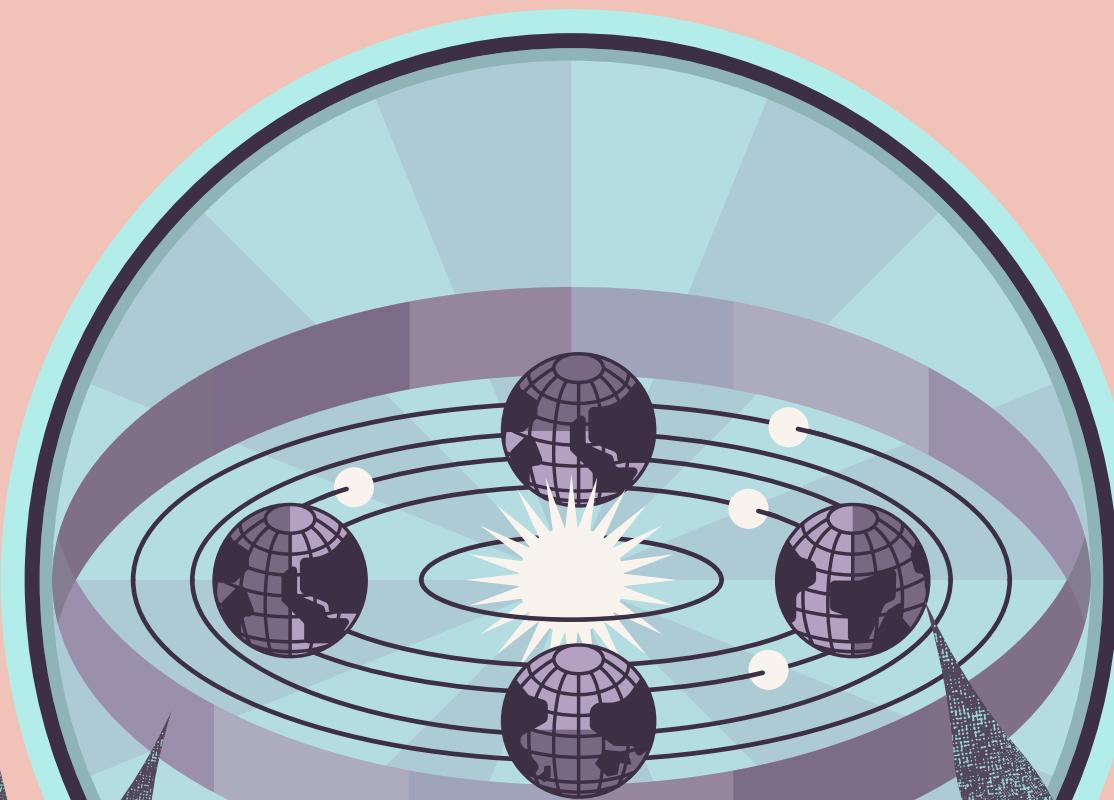
RandomForest

MODEL NAME	TRAIN SCORE	TEST SCORE	TRAIN R2	TEST R2	TRAIN MAE	TEST MAE
RIDGE	0.036	0.037	0.697	0.680	0.032	0.032
RANDOM FOREST	0.053	0.053	0.339	0.329	0.044	0.045

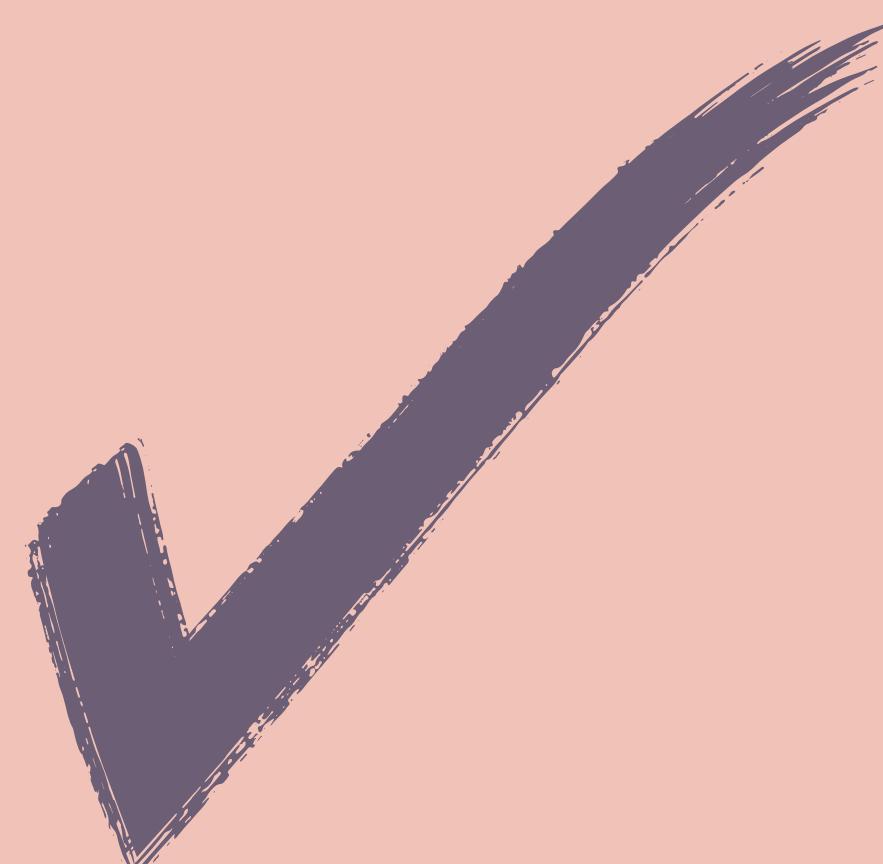


Modeling

RandomForest



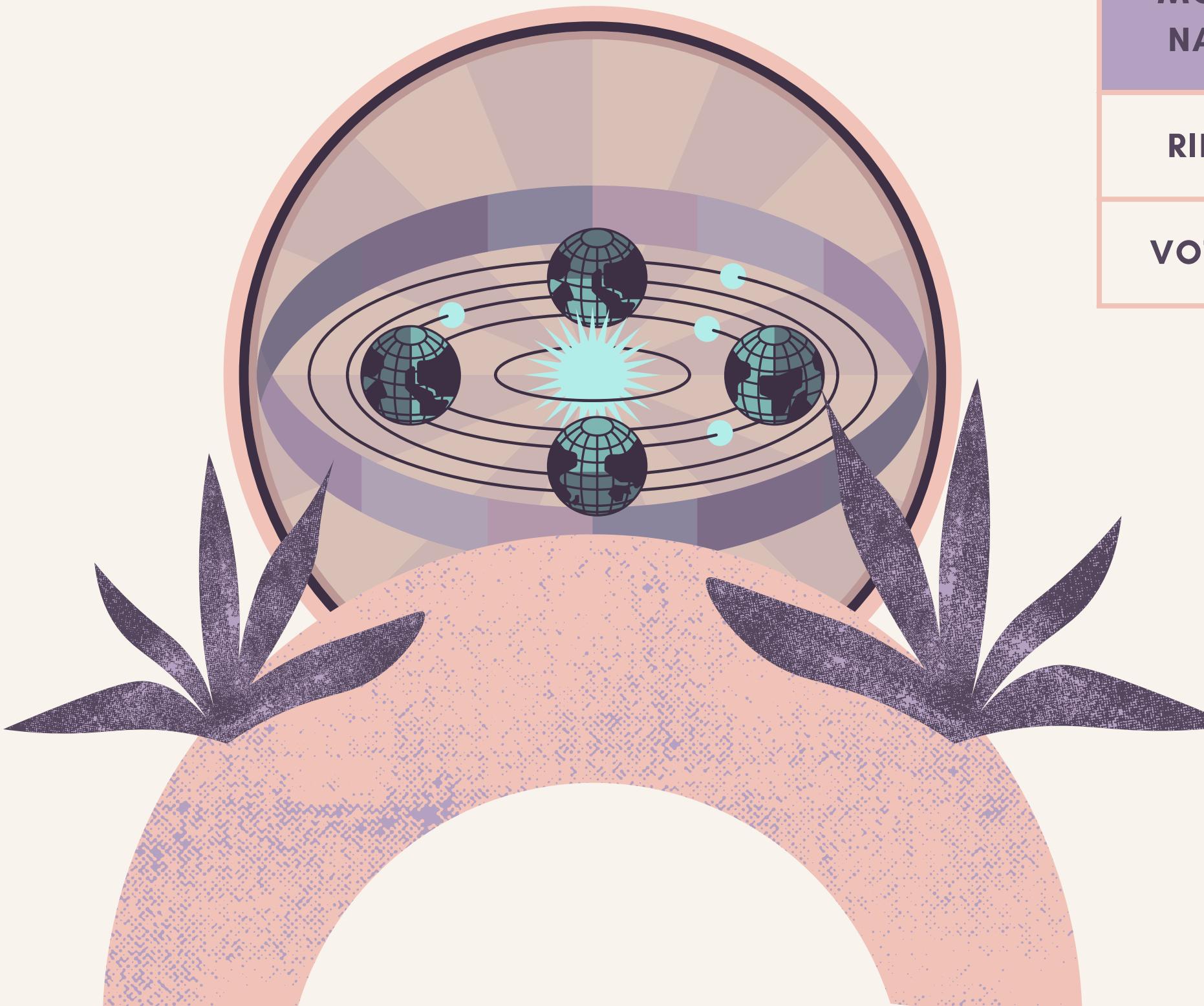
MODEL NAME	TRAIN SCORE	TEST SCORE	TRAIN R2	TEST R2	TRAIN MAE	TEST MAE
RIDGE	0.036	0.037	0.697	0.680	0.032	0.032
NEURAL NETWORK	0.045	0.045	0.700	0.700	0.032	0.032



Modeling

VotingRegressor

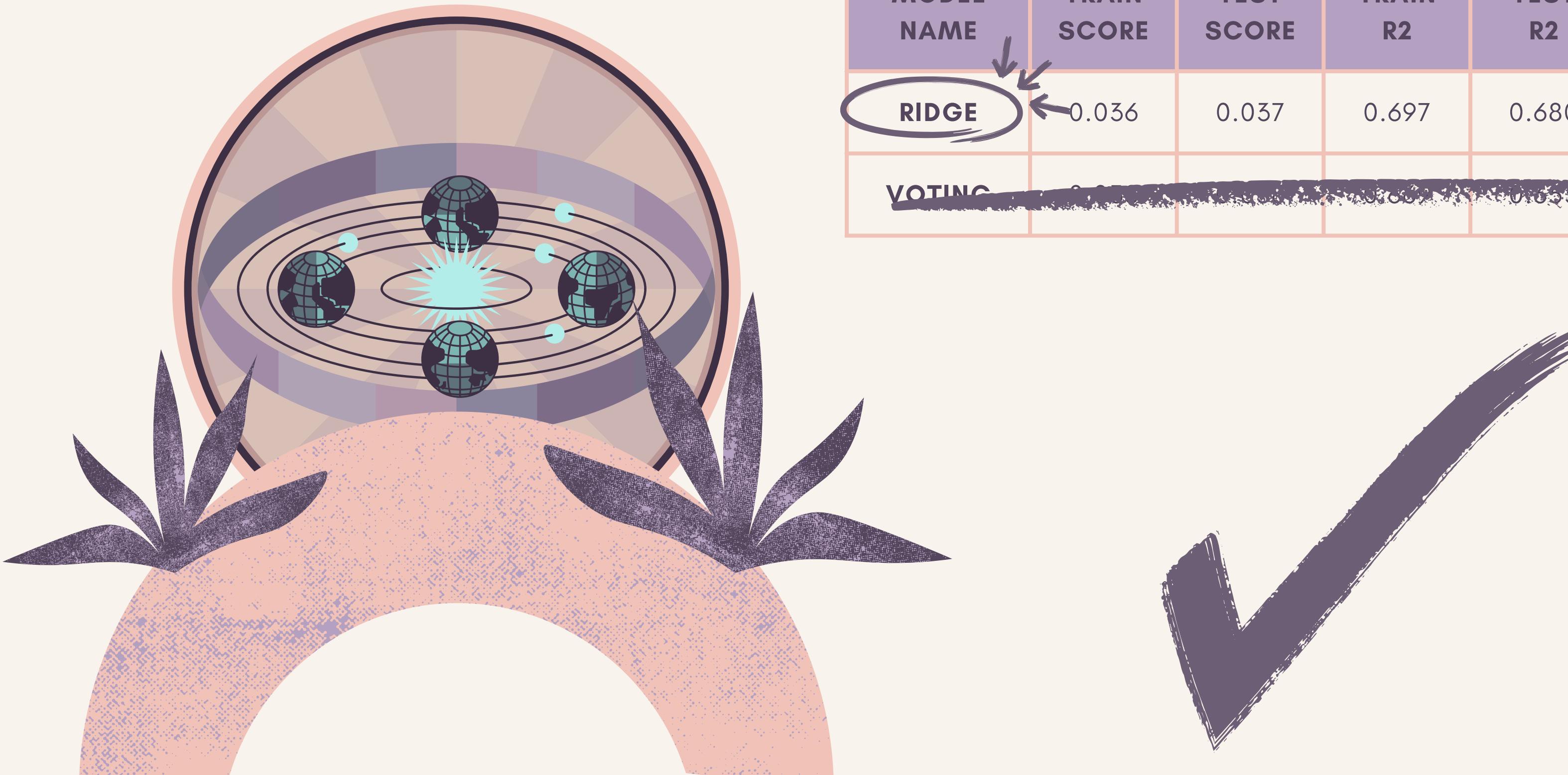
MODEL NAME	TRAIN SCORE	TEST SCORE	TRAIN R2	TEST R2	TRAIN MAE	TEST MAE
RIDGE	0.036	0.037	0.697	0.680	0.032	0.032
VOTING	0.038	0.038	0.669	0.655	0.033	0.034



Modeling

VotingRegressor

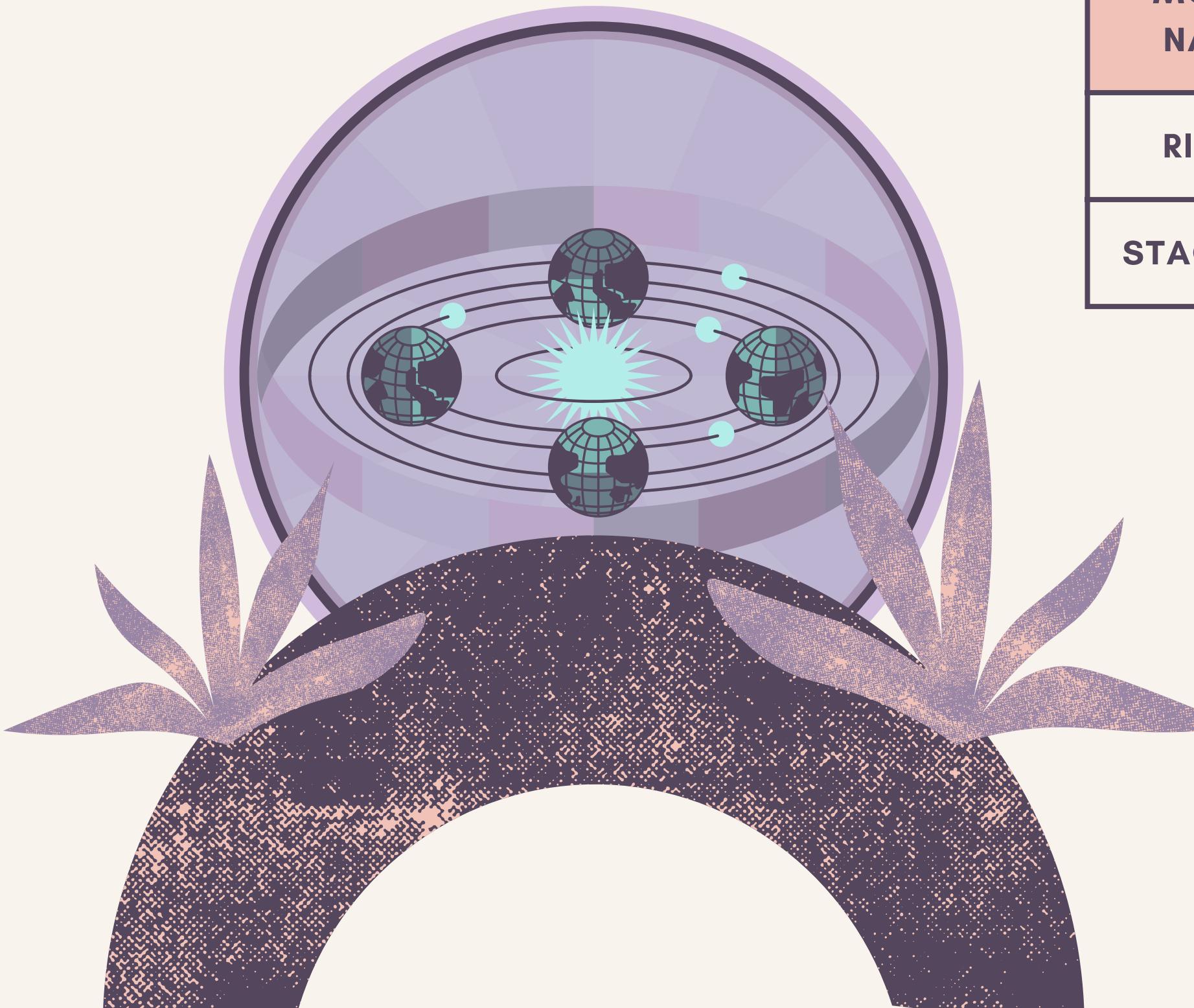
MODEL NAME	TRAIN SCORE	TEST SCORE	TRAIN R2	TEST R2	TRAIN MAE	TEST MAE
RIDGE	0.036	0.037	0.697	0.680	0.032	0.032
VOTING	0.034	0.034	0.697	0.680	0.032	0.032



Modeling

StackingRegressor

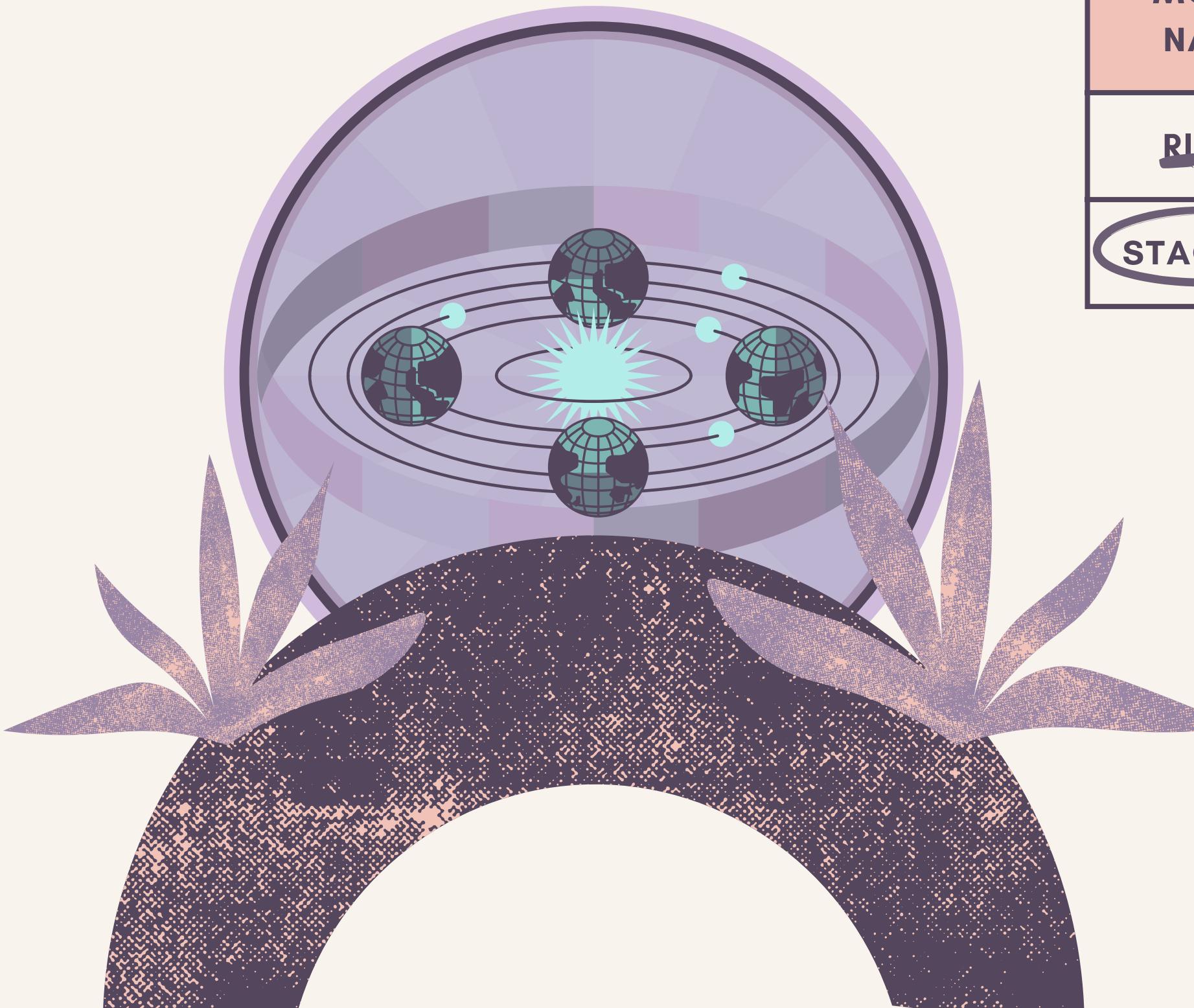
MODEL NAME	TRAIN SCORE	TEST SCORE	TRAIN R2	TEST R2	TRAIN MAE	TEST MAE
RIDGE	0.036	0.037	0.697	0.680	0.032	0.032
STACKING	0.036	0.037	0.697	0.680	0.032	0.032



Modeling

StackingRegressor

MODEL NAME	TRAIN SCORE	TEST SCORE	TRAIN R2	TEST R2	TRAIN MAE	TEST MAE
RIDGE	0.970	0.969	0.973	0.966	0.032	0.032
STACKING	0.036	0.037	0.697	0.680	0.032	0.032



Results

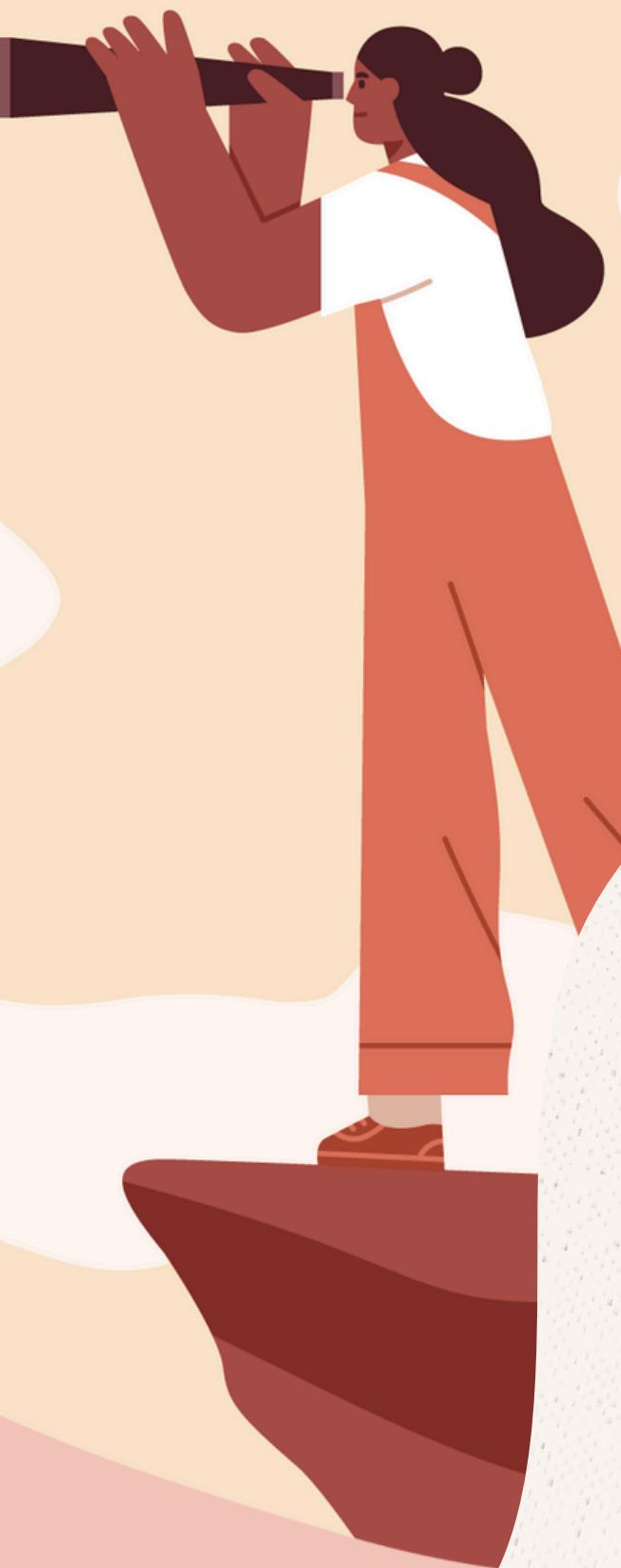
Holdout Set

Base	Final
RMSE: 1.68e+09	RMSE: 0.037
R-Squared: -6.707e+18	R-Squared: 0.698
MAE: 1.909e+07	MAE: 0.032

Looking Ahead

- 2
Deeper geolocation analysis
- 3
SDOH measures by year

- 1
Computational resources
- 4
Upscaled deployment



Deployment

Thank you!



CONTACT WITH ANY BUSINESS INQUIRIES

ELINARANKOVA@GMAIL.COM

GITHUB | BLOG | LINKEDIN