

## Step 1: Data

1. 34 columns each having 3047 data samples are included in the given dataset.

2. It will predict the target death rates.

3. Minimum and Maximum values in each respective row.

1	avgAnnCount	6.0
2	avgDeathsPerYear	3
3	TARGET_deathRate	59.7
4	incidenceRate	201.3
5	medIncome	22640
6	popEst2015	827
7	povertyPercent	3.2
8	studyPerCap	0.0
9	binmedInc	(34218.1, 37413.8]
10	MedianAge	22.3
11	MedianAgeMale	22.4
12	MedianAgeFemale	22.3
13	Geography	Abbeville County, South Carolina
14	AvgHouseholdSize	0.0221
15	PercentMarried	23.1
16	PctNoHS18_24	0.0
17	PctHS18_24	0.0
18	PctSomeCol18_24	7.1
19	PctBachDeg18_24	0.0
20	PctHS25_Over	7.5
21	PctBachDeg25_Over	2.5
22	PctEmployed16_Over	17.6
23	PctUnemployed16_Over	0.4
24	PctPrivateCoverage	22.3
25	PctPrivateCoverageAlone	15.7
26	PctEmpPrivCoverage	13.5
27	PctPublicCoverage	11.2
28	PctPublicCoverageAlone	2.6
29	PctWhite	10.199155
30	PctBlack	0.0
31	PctAsian	0.0
32	PctOtherRace	0.0
33	PctMarriedHouseholds	22.99249
34	BirthRate	0.0

Minimum Values

1	avgAnnCount	38150.0
2	avgDeathsPerYear	14010
3	TARGET_deathRate	362.8
4	incidenceRate	1206.9
5	medIncome	125635
6	popEst2015	10170292
7	povertyPercent	47.4
8	studyPerCap	9762.308998
9	binmedInc	[22640, 34218.1]
10	MedianAge	624.0
11	MedianAgeMale	64.7
12	MedianAgeFemale	65.7
13	Geography	Zavala County, Texas
14	AvgHouseholdSize	3.97
15	PercentMarried	72.5
16	PctNoHS18_24	64.1
17	PctHS18_24	72.5
18	PctSomeCol18_24	79.0
19	PctBachDeg18_24	51.8
20	PctHS25_Over	54.8
21	PctBachDeg25_Over	42.2
22	PctEmployed16_Over	80.1
23	PctUnemployed16_Over	29.4
24	PctPrivateCoverage	92.3
25	PctPrivateCoverageAlone	78.9
26	PctEmpPrivCoverage	70.7
27	PctPublicCoverage	65.1
28	PctPublicCoverageAlone	46.6
29	PctWhite	100.0
30	PctBlack	85.947799
31	PctAsian	42.619425
32	PctOtherRace	41.930251
33	PctMarriedHouseholds	78.075397
34	BirthRate	21.326165

Maximum Values

The minimum value in the entire dataset = 0.0

The maximum value in the entire dataset = 10170292.0

4. There are 34 columns, hence 34 features in each data sample.

5. Columns with missing values along with the total cells missing:

PctSomeCol18\_24 → 2285, PctEmployed16\_Over → 152, PctPrivateCoverageAlone → 609

6. The label of this dataset is "TARGET\_deathRate."

7. 70% for training, 15% for validation, and 15% for testing.

8. I began by reading a CSV file called 'cancer\_reg.csv' and stored its contents in a DataFrame named **df**.

Next, I needed to address missing values in the dataset. To do this, I calculated the sum of missing values for each row in **df** using **df.isnull().sum(axis=1)**. Then, I looped through the columns, a total of 34, to check for any missing values in each column. Whenever I found a column with missing values, I filled those gaps by replacing them with the mean value of that column using **df.iloc[:, i].fillna(df.iloc[:, i].mean(), inplace=True)**.

Following this, I performed label encoding on two specific columns, 'Geography' and 'binnedInc', using Scikit-Learn's **LabelEncoder**. This transformation converted the categorical data in those columns into numerical format.

To define the features (X) and target (y) for my machine learning model, I excluded the 'TARGET\_deathRate' column from the DataFrame to create X. I set y to be the 'TARGET\_deathRate' column.

Then, I applied a logarithmic transformation to both X and y using **np.log1p()**. This transformation helped me make the data more suitable for modeling, especially if the data had skewed distributions or large ranges.

For evaluating the model's performance, I split the data into three sets: a training set (70% of the data), a validation set (15%), and a testing set (15%). This was accomplished with the **train\_test\_split** function, ensuring that I had separate datasets for training and testing.

Standardization was the next step in my data preparation. I used **StandardScaler** to scale the feature data. First, I fitted the scaler to the training data, and then I applied the same scaler to transform the validation and testing data. Standardization is important for ensuring that all features have zero mean and unit variance, which can aid certain machine learning algorithms in performing better.

Lastly, I performed clipping on the predictions and validation data to prevent them from going below a specified minimum value. This was done using **np.clip()**, which constrained the values within a predefined range.

## Step 2: Model

Model	MSE
Linear regression	0.0046
ANN-oneL-16	0.0049
ANN-twoL-32-8	0.0037
ANN-threeL-32-16-8	0.0033
ANN-fourL-32-16-8-4	0.0029

## Step 3: Objective

Mean Squared Error (MSE) is the loss function I used to train my models.

## Step 4: Optimization

I did some research about Adam optimizer on my own and hence have used Adam optimizer for my model to train which also helped me to increase my accuracy.

## Step 5: Model Selection

Model	LR: 0.1 (100 epochs, 64 batch size)	LR: 0.01 (1000 epochs, 64 batch size)	LR: 0.001 (1000 epochs, 64 batch size)	LR: 0.0001 (1000 epochs, 64 batch size)
ANN-oneL-16	0.76	0.84	0.79	0.50
ANN-twoL-32-8	0.79	0.84	0.72	0.51
ANN-threeL-32-16-8	0.81	0.88	0.72	0.54
ANN-fourL-32-16-8-4	0.82	0.88	0.73	0.69

The best performing model was the ANN-fourL-32-16-8-4 model at 1000 epochs, 64 batch size and 0.01 learning rate which gave me a R-squared score of 0.88.

## Step 6: Model Performance

For the R-squared score table (Step 5):

1. Higher learning rates (e.g., LR: 0.1) generally lead to better R-squared scores, but there's a trade-off with stability and convergence.
2. Increasing the number of epochs (e.g., 1000 epochs) tends to improve model performance, as reflected in the higher R-squared scores.
3. The choice of model architecture (e.g., the number of hidden layers and units) also affects performance, with more complex models (e.g., ANN-fourL-32-16-8-4) generally outperforming simpler ones (e.g., ANN-oneL-16).
4. Very low learning rates (e.g., LR: 0.0001) can lead to suboptimal model performance, as seen in the lower R-squared scores.

For the MSE table (Step 2):

1. The ANN models (Artificial Neural Networks) are performing better than the simple Linear Regression model in terms of MSE. This suggests that the neural networks are capturing more complex relationships within the data.
2. As you increase the depth of the neural network (going from one hidden layer to four hidden layers), the MSE decreases, indicating that deeper architectures can model the data more effectively.
3. Among the neural network models, "ANN-fourL-32-16-8-4" has the lowest MSE, suggesting that it is the best-performing model among the ones listed.
4. The specific choice of model and its architecture can have a significant impact on the quality of predictions, and it's important to select the model that best fits your dataset and problem.