## **Data Collection and Preprocessing**

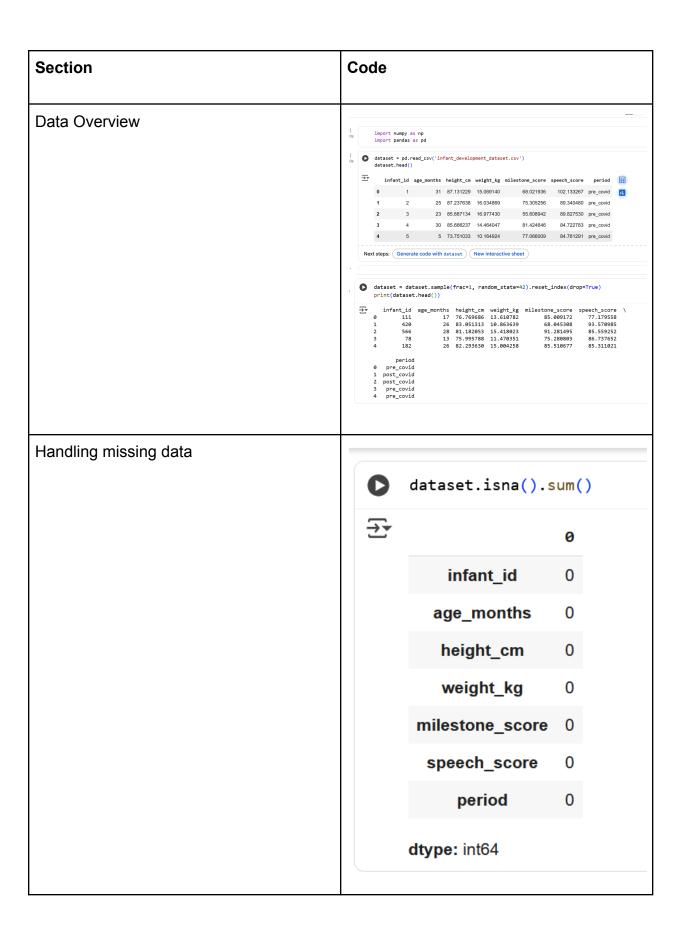
**Project Name :** Covid - 19 Infant Growth Analysis and Prediction

Prepare the infant development dataset for machine learning classification using TabPFN.

## PREPROCESSING STEPS:

| SECTION                | DESCRIPTION  |
|------------------------|--|
| Data Overview          | Loaded the dataset, shuffled dataset, inspected columns, and checked missing values.   |
| Handling missing data  | Filled missing numeric values (age_months, height_cm, weight_kg, speech_score, milestone_score) with their column means. For categorical column period replace missing values with the mode. |
| Feature & Target Split | Split dataset into features x and y  |
| Encoding               | Applied LabelEncoder to transform categorical target labels into numeric form.   |
| Splitting dataset      | Divided data into training (75%) and testing (25%) sets  |
| Model preparation      | Installed and initialized TabPFNClassifier for training.   |
| Visualization          | Generated scatter plots to compare training and testing predictions versus actual labels.  |

## **Data Preprocessing Code Screenshots:**



```
Feature & Target Split
                                                                                                                                                                                                                                     x = dataset.iloc[:, :-1].values
                                                                                                                                                                                                                                     → array([[111.
                                                                                                                                                                                                                                                                             11. , 17. ,
85.00917188, 77.17955804],
                                                                                                                                                                                                                                                                                                                                                            76.76968601, 13.61078245,
                                                                                                                                                                                                                                                                            420. , 26. ,
68.04530769, 93.57098541],
                                                                                                                                                                                                                                                                                                                                                              83.05131333, 10.86363906,
                                                                                                                                                                                                                                                                       [420.
                                                                                                                                                                                                                                                                           566. , 28. ,
91.28149546, 85.55925192],
                                                                                                                                                                                                                                                                                                                                                              81.18205336, 15.41802335,
                                                                                                                                                                                                                                                                           271. , 14. , 76.82129208, 73.9939117 ], 436. , 12. ,
                                                                                                                                                                                                                                                                                                                                                              72.98423838, 12.75260549,
                                                                                                                                                                                                                                                                                                                                                            74.07420025, 8.29130125,
                                                                                                                                                                                                                                                                            49.37666341, 76.08972237],
                                                                                                                                                                                                                                                                            49.57666541, 76.68972257],
103. , 17. , 8
87.68727492, 78.75725317]])
                                                                                                                                                                                                                                                                                                                                                                80.89907367, 12.23801102,
                                                                                                                                                                                                                                                                       [103.
                                                                                                                                                                                                                                   y = dataset.iloc[:, -1].values
                                                                                                                                                                                                                                                              'post_covid', 'pre_covid', 'pre_covid', 'during_covid',
'post_covid', 'pre_covid', 'during_covid',
'during_covid', 'pre_covid', 'pre_covid',
'during_covid', 'during_covid', 'pre_covid',
'post_covid', 'post_covid', 'post_covid', 'post_covid',
'post_covid', 'pre_covid', 'post_covid', 'during_covid',
'during_covid', 'pre_covid', 'post_covid', 'during_covid',
'post_covid', 'during_covid', 'post_covid', 'pre_covid',
'post_covid', 'pre_covid', 'post_covid', 'pre_covid',
'post_covid', 'pre_covid', 'pre_covid', 'pre_covid',
'post_covid', 'pre_covid', 'pre_covid', 'pre_covid',
'pre_covid', 'during_covid', 'during_covid', 'pre_covid',
'pre_covid', 'during_covid', 'during_covid', 'pre_covid',
'pre_covid', 'during_covid', 'during_covid', 'during_covid',
'pre_covid', 'pre_covid', 'pre_covid', 'during_covid',
'during_covid', 'during_covid', 'during_covid',
'during_covid', 'pre_covid', 'pre_covid', 'pre_covid',
'pre_covid', 'pre_covid', 'pre_covid', 'pre_covid',
'pre_covid', 'pre_covid', 'pre_covid', 'pre_covid',
'pre_covid', 'pre_covid', 'pre_covid', 'pre_covid',
'during_covid', 'pre_covid', 'pre_covid', 'pre_covid',
'during_covid', 'pre_covid', 'pre_covid', 'pre_covid',
'during_covid', 'pre_covid', 'pre_covid', 'pre_covid',
'during_covid', 'pre_covid', 'pre_covid', 'pre_covid', 'pre_covid',
'during_covid', 'pre_covid', 'pr
                                                                                                                                                                                                                                   ∓÷
Encoding
                                                                                                                                                                                                                                           #encoding missing values
                                                                                                                                                                                                                                                         dataset.fillna({
                                                                                                                                                                                                                                                                       'age_months': dataset['age_months'].mean(),
                                                                                                                                                                                                                                                                       'height_cm': dataset['height_cm'].mean(),
                                                                                                                                                                                                                                                                      'weight_kg': dataset['weight_kg'].mean(),
                                                                                                                                                                                                                                                                      'speech_score': dataset['speech_score'].mean(),
                                                                                                                                                                                                                                                                      'period': dataset['period'].mode()[0] }, inplace=True)
                                                                                                                                                                                                                                     # Encoding categorical data
                                                                                                                                                                                                                                                 from sklearn.preprocessing import LabelEncoder
                                                                                                                                                                                                                                                 le = LabelEncoder()
                                                                                                                                                                                                                                                 y = le.fit_transform(y)
```

