

Import Libraries we want to use

```
[2]: import pandas as pd
#For virtualization
import seaborn as sns
import matplotlib.pyplot as plt
```

Loading Datasets

```
[4]: # Load the baby weights dataset
dataset = pd.read_csv('baby-weights-dataset.csv')
dataset
```

	ID	SEX	MARITAL	FAGE	GAINED	VISITS	MAGE	FEDUC	MEDUC	TOTALP	...	HYPERC	HYPERRP	ECLAMP	CERVIX	PINFANT	PRETERM	RENAL	RHS
0	2001	2	1	33	26.0	10	34	12.0	4	2	—	0	0	0	0	0	0	0	0
1	2002	2	2	19	40.0	10	18	11.0	12	1	—	0	0	0	0	0	0	0	0
2	2003	2	1	33	16.0	14	31	16.0	16	2	—	0	0	0	0	0	0	0	0
3	2004	1	1	25	40.0	15	28	12.0	12	3	—	0	0	0	0	0	0	0	0
4	2005	1	2	21	60.0	13	20	12.0	14	2	—	0	1	0	0	0	0	0	0
...	—	—	—	—	—	—	—	—	—
101395	103396	1	2	36	0.0	9	34	3.0	12	4	—	0	0	0	0	0	0	0	0
101396	103397	2	2	21	39.0	11	19	12.0	9	2	—	0	0	0	0	0	0	0	0
101397	103398	2	1	27	37.0	15	22	12.0	12	2	—	0	0	0	0	0	0	0	0
101398	103399	1	1	27	26.0	12	24	12.0	14	1	—	0	0	0	0	0	0	0	0
101399	103400	1	2	20	31.0	15	17	12.0	11	1	—	0	0	0	0	0	0	0	0

101400 rows × 37 columns

Data Preprocessing**Data Exploration**

describing columns in dataset .. provides a summary of each feature, including mean, standard deviation, and quartiles.

```
[7]: dataset.describe()
```

	ID	SEX	MARITAL	FAGE	GAINED	VISITS	MAGE	FEDUC	MEDUC	TOTALP	...	HYP
count	101400.000000	101400.000000	101400.000000	101399.000000	101400.000000	101399.000000	101400.000000	101400.000000	101400.000000	101400.000000	...	101400
mean	52700.500000	1.485671	1.303817	30.174477	30.283040	12.436943	27.736312	12.926883	13.256489	2.378462	...	0
std	29271.802985	0.500349	0.459907	6.775576	13.615468	3.728901	5.957369	2.926582	2.932693	1.490272	...	0
min	2001.000000	1.000000	1.000000	14.000000	0.000000	0.000000	11.000000	0.000000	0.000000	1.000000	...	0
25%	27350.750000	1.000000	1.000000	25.000000	21.000000	10.000000	23.000000	12.000000	12.000000	12.000000	...	0
50%	52700.500000	1.000000	1.000000	30.000000	30.000000	12.000000	28.000000	12.000000	13.000000	13.000000	...	0
75%	78050.250000	2.000000	2.000000	35.000000	39.000000	15.000000	32.000000	16.000000	16.000000	16.000000	...	0
max	103400.000000	9.000000	2.000000	74.000000	98.000000	49.000000	53.000000	17.000000	17.000000	20.000000	...	1

8 rows × 35 columns

dataset size : number of columns , rows

```
[9]: dataset.shape
```

[9]: (101400, 37)

dataset information

```
[11]: dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 101400 entries, 0 to 101399
Data columns (total 37 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   ID          101400 non-null  int64  
 1   SEX         101400 non-null  int64  
 2   MARITAL     101400 non-null  int64  
 3   FAGE        101400 non-null  int64  
 4   GAINED      101399 non-null  float64 
 5   VISITS     101400 non-null  int64  
 6   MAGE        101400 non-null  int64  
 7   FEDUC       101399 non-null  float64 
 8   MEDUC       101400 non-null  int64  
 9   TOTALP      101400 non-null  int64  
 10  BDEAD       101400 non-null  int64  
 11  TERMS       101400 non-null  int64  
 12  LOUTCOME    101400 non-null  int64  
 13  WEEKS       101399 non-null  float64 
 14  RACEMOM     101400 non-null  int64  
 15  RACEDAD     101400 non-null  int64  
 16  HISPMON     101400 non-null  object 
 17  HISPDAD     101400 non-null  object 
 18  CIGNUM      101399 non-null  float64 
 19  DRINKNM     101400 non-null  int64  
 20  ANEMIA      101400 non-null  int64  
 21  CARDIAC     101400 non-null  int64  
 22  ACLUNG      101400 non-null  int64  
 23  DIABETES    101400 non-null  int64  
 24  HERPES      101400 non-null  int64  
 25  HYDRAM      101399 non-null  float64 
 26  HEMOGLOB    101400 non-null  int64  
 27  HYPERC      101400 non-null  int64  
 28  HYPERRP     101400 non-null  int64  
 29  ECLAMP      101400 non-null  int64  
 30  CERVIX      101400 non-null  int64  
 31  PINFANT     101400 non-null  int64  
 32  PRETERM     101400 non-null  int64  
 33  RENAL       101400 non-null  int64  
 34  RHSEN       101400 non-null  int64  
 35  UTERINE     101400 non-null  int64  
 36  BWWEIGHT    101400 non-null  float64 
dtypes: float64(6), int64(29), object(2)
memory usage: 28.6+ MB
```

```
[12]: # make column name Lowercase
dataset.columns = dataset.columns.str.lower()
dataset.columns
```

```
[12]: Index(['id', 'sex', 'marital', 'fage', 'gained', 'visits', 'mage', 'feduc',
       'meduc', 'totalp', 'bdead', 'terms', 'loutcome', 'weeks', 'racemom',
       'racedad', 'hispmom', 'hispdad', 'cignum', 'drinknm', 'anemia',
       'cardiac', 'aclung', 'diabetes', 'herpes', 'hydram', 'hemoglob',
       'hyperch', 'hyperpr', 'eclamp', 'cervix', 'pinfant', 'preterm', 'renal',
       'rhsen', 'uterine', 'bwweight'],
      dtype='object')
```

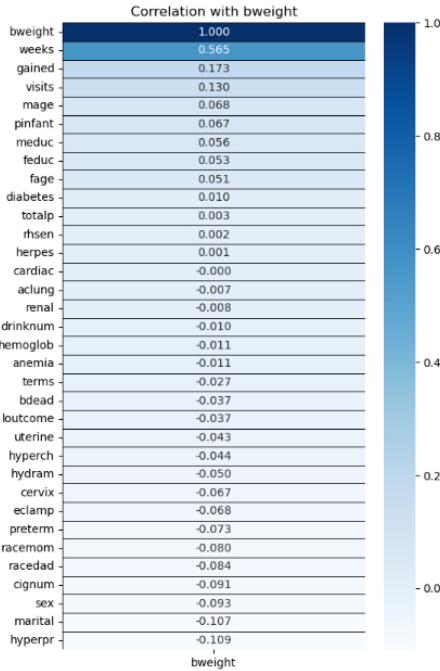
```
[13]: # Drop the id column
```

```
[14]: dataset.drop(columns=['id'], inplace=True) # can use index nested
# numeric columns
numeric_columns = dataset.select_dtypes(include=['number']).columns

# calc correlation
correlation_matrix = dataset[numeric_columns].corr()
correlation_with_bweight = correlation_matrix[['bweight']].sort_values(by='bweight', ascending=False)

# heatmap
plt.figure(figsize=(6, 10))
sns.heatmap(correlation_with_bweight, annot=True, cmap='Blues', linewidths=0.5, linecolor='black', fmt=".3f")
plt.title('Correlation with bweight')
plt.yticks(rotation=0)

# Save the heatmap
plt.savefig('heatmap_bweight.png')
```



Checking the distribution of data in categorical different classes

```
[16]: # Check the distribution of the 'sex' column
dataset['sex'].value_counts()

[16]: sex
1    52160
2    49239
9     1
Name: count, dtype: int64
we notice outlier in value 9

[18]: # Check the distribution of the 'marital' column
dataset['marital'].value_counts()

[18]: marital
1    76593
2    38887
Name: count, dtype: int64

[19]: # Check the distribution of the 'racemom' column
dataset['racemom'].value_counts()

[19]: racemom
1    78343
2    18015
8     3260
3     1196
6     224
7      159
4      132
5      42
6      28
Name: count, dtype: int64

[20]: # Value counts categorical columns
for col in ['racedad', 'hispmom', 'hispdad']:
    print(dataset[col].value_counts())
    print("-" * 20)

racedad
1    76285
2    28500
8     2915
3     1190
0     285
4     123
7      97
9      50
6      18
5      17
Name: count, dtype: int64
-----
hispmom
N    85040
M    11290
S     3757
P     981
C     185
O     135
U      92
Name: count, dtype: int64
-----
hispdad
N    84301
M    12039
S     3771
P     884
C     186
O     118
U     181
Name: count, dtype: int64
-----
```

Checking For missing Values

```
[22]: dataset.isnull().sum() #sum of null values

[22]: sex          0
       marital      0
       fage          0
       gained        1
       .....
```

```

vars   v
mage  0
feduc 0
meduc 0
totalp 0
bdead 0
terms 0
loutcome 0
weeks 1
racemom 0
racedad 0
hispmom 0
hispdad 0
cignum 1
drinknum 0
anemia 0
cardiac 0
aclung 0
diabetes 0
herpes 0
hydram 1
hemoglobin 0
hyperch 0
hyperpr 0
clamp 0
cervix 0
pinfant 0
preterm 0
renal 0
rhens 0
uterine 0
bwheight 0
dtype: int64
[23]: dataset.isnull().sum().sum() # num of rows has null values
[23]: 5

```

Checking for duplication

```

[25]: dataset.duplicated().sum() # sum of duplicated values
[25]: 2
[26]: dataset[dataset.duplicated()]
[26]:
   sex marital fage gained visits mage feduc meduc totalp bdead ... hyperch hyperpr clamp cervix pinfant preterm renal rhens uterine bwheight
45040 2     1    20   25.0    12    22   12.0    14    2    0   ...      0      0      0      0      0      0      0      0      0      0      0      0      0      0      1.4375
54306 1     2    28   11.0    5    26   10.0    10    6    0   ...      0      0      0      0      0      0      0      0      0      0      0      0      0      0      3.6250
2 rows x 36 columns

```

Data Cleaning

Handling missing values and duplicates ensures that the dataset is free from errors and inconsistencies, which is crucial for accurate modeling.

Handle Missing values

```

[28]: # Drop rows with any missing values
dataset.dropna(inplace=True)
dataset.isnull().sum() # should be no missing values
[28]:
   sex     0
marital  0
fage     0
gained   0
visits   0
mage     0
feduc    0
meduc    0
totalp   0
bdead    0
terms    0
loutcome 0
weeks   0
racemom  0
racedad  0
hispmom  0
hispdad  0
cignum   0
drinknum 0
anemia   0
cardiac  0
aclung   0
diabetes 0
herpes   0
hydram   0
hemoglobin 0
hyperch  0
hyperpr  0
clamp   0
cervix  0
pinfant 0
preterm  0
renal   0
rhens   0
uterine 0
bwheight 0
dtype: int64

```

Handle duplicated values with drop/delete

```

[30]: print(len(dataset))
dataset = dataset.drop_duplicates(keep='first')
print(len(dataset))
101395
101393

```

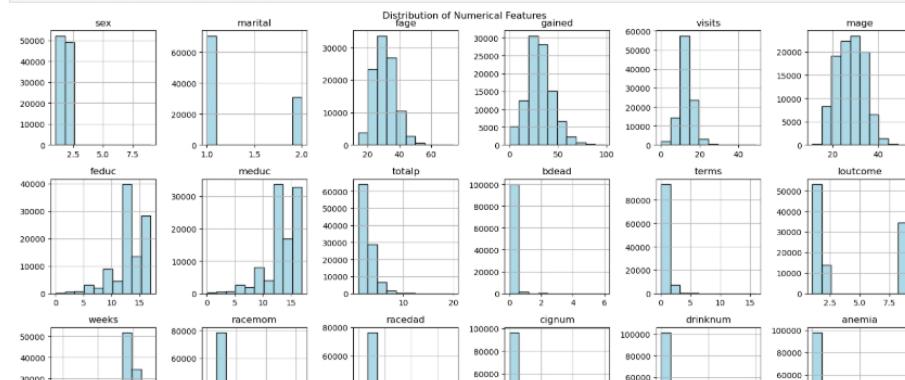
Data Visualization

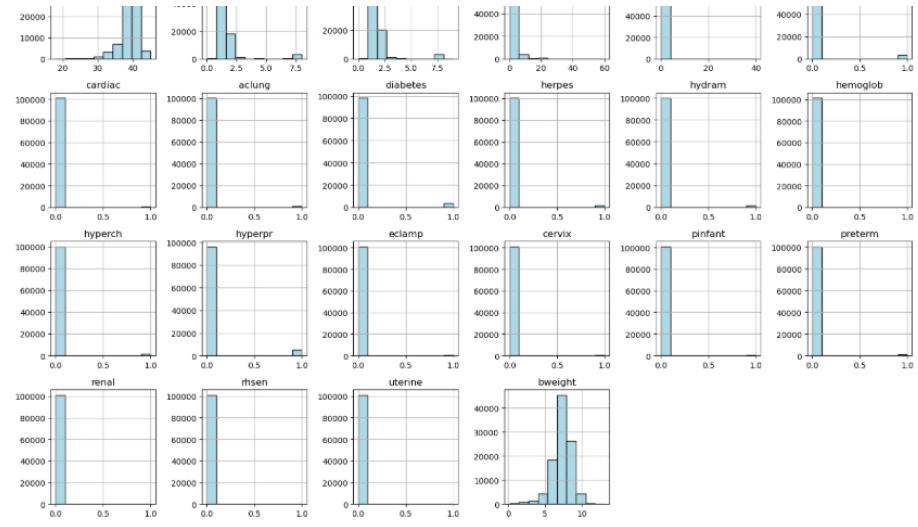
show distribution of data

```

[33]: # Plot histograms for numerical features
numeric_columns = dataset.select_dtypes(include=['float64', 'int64']).columns
dataset[numeric_columns].hist(bins=10, figsize=(16, 16), color='lightblue', edgecolor='black')
plt.suptitle('Distribution of Numerical Features')
plt.tight_layout()
plt.savefig('hist_before_handel_outlier.png')

```





Notice outlier & unusual distribution on some attributes like totalap , fage , visits , sex and so on , each histogram explain data distribution and if exist very large or very small data

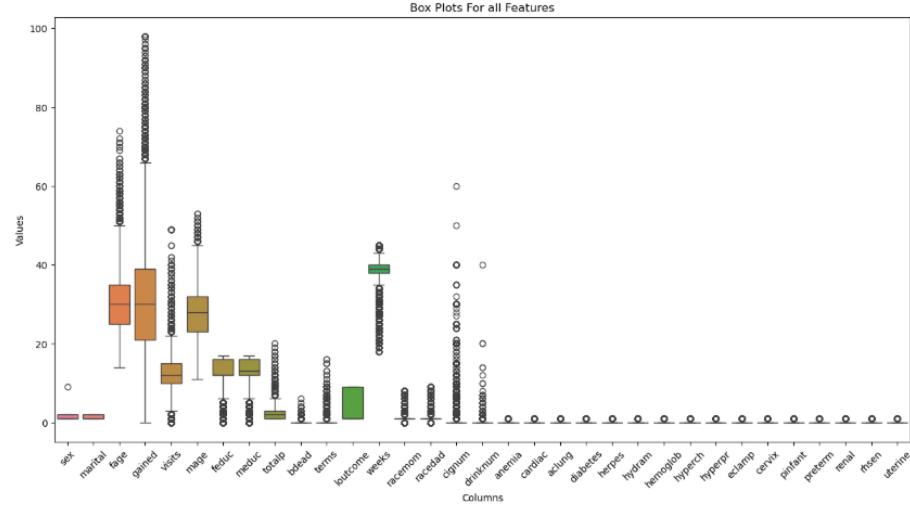
Handle outlier values

use sns to make boxplot and show data outliers

```
[37]: # boxplot for all columns in the dataset
plt.figure(figsize=(16, 8))
sns.boxplot(data=dataframe.iloc[:, :-1]) # except last/object column
plt.title('Box Plots For all Features')
plt.xticks(rotation=45)
plt.xlabel('Values')
plt.ylabel('Columns')

# plt.savefig('single_boxplot.png')
```

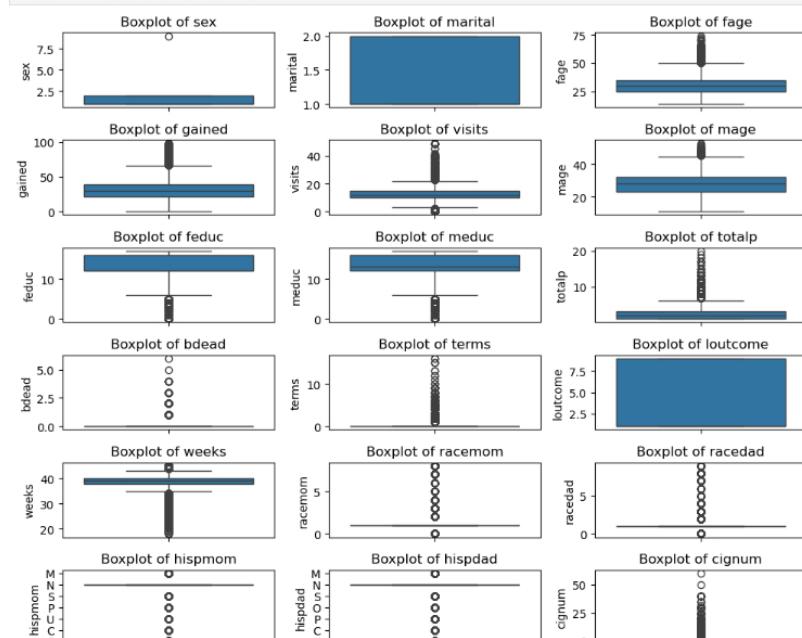
[37]: Text(0.5, 0, 'Columns')

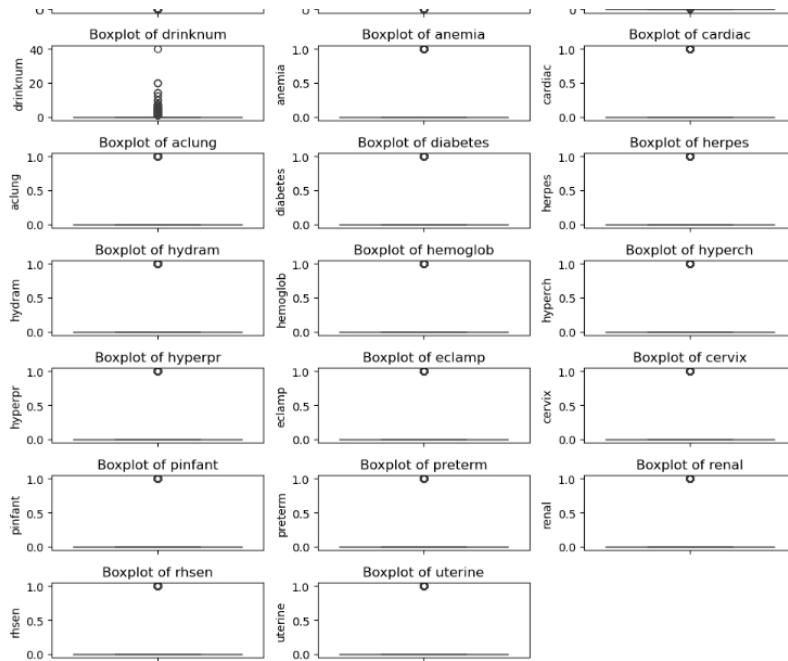


```
[38]: plt.figure(figsize=(10, 16))

# boxplots for each column
for i, column in enumerate(dataframe.columns[:-1], 1): # except the last/object column
    plt.subplot(12, 3, i) # in 12 rows , 3 columns
    sns.boxplot(data=dataframe, y=column)
    plt.title(f'Boxplot of {column}')
plt.tight_layout() # spacing between subplots

# plt.savefig('individual_boxplots.png')
```





```
[39]: # handle sex column has outlier with value 9
print("len of dataset before drop sex column outlier : ", len(dataset))
dataset = dataset[dataset['sex'] != 9]
print("len of dataset after drop sex column outlier : ", len(dataset))

len of dataset before drop sex column outlier : 101393
len of dataset after drop sex column outlier : 101392
```

```
[40]: def handle_binary_columns(dataset, binary_columns):
    # handle binary columns - only rows with 0 or 1 values
    for col in binary_columns:
        dataset = dataset[dataset[col].isin([0, 1])]
```

```
return dataset
```

```
[41]: print(len(dataset))
dataset = handle_binary_columns(dataset , ['anemia', 'cardiac', 'aclung', 'diabetes', 'herpes']) # verify but boxplot display no outlier
print(len(dataset))
```

```
101392
```

```
101392
```

```
[42]: def remove_outliers(dataset, columns):
    for col in columns:
        if dataset[col].dtype != 'object':
            q1_value = dataset[col].quantile(0.25)
            q3_value = dataset[col].quantile(0.75)
            IQR = q3_value - q1_value
            max_value = q3_value + (1.5 * IQR)
            min_value = q1_value - (1.5 * IQR)
            dataset = dataset[(dataset[col] >= min_value) & (dataset[col] <= max_value)]
```

```
[43]: print(len(dataset))
dataset = remove_outliers(dataset , ['fage', 'gained', 'visits', 'mage', 'meduc', 'totalp'])
print(len(dataset))
```

```
101392
```

```
94954
```

```
[44]: dataset[['drinknum', 'cignum', 'hispmom', 'hispdad', 'racemom', 'racedad', 'weeks']].describe()
```

	drinknum	cignum	hispmom	hispdad	racemom	racedad	weeks
count	94054.000000	94054.000000	94054.000000	94054.000000	94054.000000	94054.000000	94054.000000
mean	0.004827	0.708433	1.441725	1.441108	38.766794		
std	0.161786	2.923391	1.302274	1.246343	2.457950		
min	0.000000	0.000000	0.000000	0.000000	18.000000		
25%	0.000000	0.000000	1.000000	1.000000	38.000000		
50%	0.000000	0.000000	1.000000	1.000000	39.000000		
75%	0.000000	0.000000	1.000000	1.000000	40.000000		
max	92.000000	60.000000	8.000000	9.000000	45.000000		

Handle outlier on 'drinknum', 'cignum' with mean value

```
[45]: # IQR for drinknum
Q1_drinknum = dataset['drinknum'].quantile(0.25)
Q3_drinknum = dataset['drinknum'].quantile(0.75)
IQR_drinknum = Q3_drinknum - Q1_drinknum
min_drinknum = Q1_drinknum - 1.5 * IQR_drinknum
max_drinknum = Q3_drinknum + 1.5 * IQR_drinknum

# mean value of drinknum
mean_drinknum = dataset['drinknum'].mean()

# mean value
dataset.loc[(dataset['drinknum'] < min_drinknum) | (dataset['drinknum'] > max_drinknum), 'drinknum'] = mean_drinknum

dataset['drinknum'].describe()
```

C:\Users\HP\AppData\Local\Temp\ipykernel_14296\302503016.py:12: FutureWarning: Setting an item of incompatible dtype is deprecated and will raise an error in a future version of pandas. Value '0.0048270142683090045' has dtype incompatible with int64, please explicitly cast to a compatible dtype first.
dataset.loc[(dataset['drinknum'] < min_drinknum) | (dataset['drinknum'] > max_drinknum), 'drinknum'] = mean_drinknum

	count	mean	std	min	25%	50%	75%	max	Name: drinknum, dtype: float64
count	94054.000000	0.004827	0.161786	0.000000	0.000000	0.000000	0.000000	92.000000	
mean	0.004827	0.004827	0.000228	0.000000	0.000000	0.000000	0.000000	0.000000	
std	0.161786	0.161786	0.000228	0.000000	0.000000	0.000000	0.000000	0.000000	
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	
25%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	
50%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	
75%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	
max	92.000000	0.004827	0.161786	0.000000	0.000000	0.000000	0.000000	92.000000	

```
[47]: # IQR for cignum
Q1_cignum = dataset['cignum'].quantile(0.25)
Q3_cignum = dataset['cignum'].quantile(0.75)
IQR_cignum = Q3_cignum - Q1_cignum
max_cignum = Q1_cignum - 1.5 * IQR_cignum
min_cignum = Q3_cignum + 1.5 * IQR_cignum

# mean value of cignum
mean_cignum = dataset['cignum'].mean()

# mean value nested of outlier
dataset.loc[(dataset['cignum'] < max_cignum) | (dataset['cignum'] > min_cignum), 'cignum'] = mean_cignum
```

```
[47]: # Verify the changes
dataset['cignum'].describe()

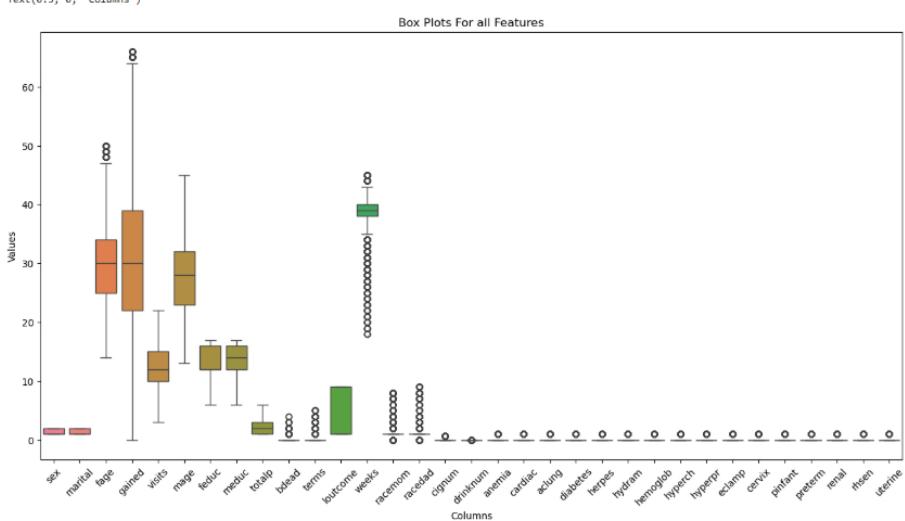
[48]: dataset[['drinknum', 'cignum', 'hispmom', 'hispdad', 'racemom', 'racedad', 'weeks']].describe()

[48]:
   drinknum    cignum    racemom    racedad    weeks
count 94054.000000 94054.000000 94054.000000 94054.000000 94054.000000
mean  0.000011  0.056657  1.441725  1.441108 38.766794
std   0.000228  0.192167 1.302274  1.246343 2.457950
min   0.000000  0.000000  0.000000  0.000000 18.000000
25%   0.000000  0.000000  1.000000  1.000000 38.000000
50%   0.000000  0.000000  1.000000  1.000000 39.000000
75%   0.000000  0.000000  1.000000  1.000000 40.000000
max   0.004827  0.708433  8.000000  9.000000 45.000000
```

```
[49]: # boxplot for all columns in the dataset
plt.figure(figsize=(16, 8))
sns.boxplot(data=dataset.iloc[:, :-1]) # except last/object column
plt.title('Box Plots For all Features')
plt.xticks(rotation=45)
plt.ylabel('Values')
plt.xlabel('Columns')

# plt.savefig('single_boxplot_after_drop_outlier.png')
```

[49]: Text(0.5, 0, 'Columns')



```
[50]: # other function to handle outlier : Clip
def clip_outliers(dataset, columns):
    for col in columns:
        if dataset[col].dtype != 'object':
            q1_value = dataset[col].quantile(0.25)
            q3_value = dataset[col].quantile(0.75)
            IQR = q3_value - q1_value
            max = q3_value + (1.5 * IQR)
            min = q1_value - (1.5 * IQR)
            dataset[col] = dataset[col].clip(lower=min, upper=max)
    return dataset
```

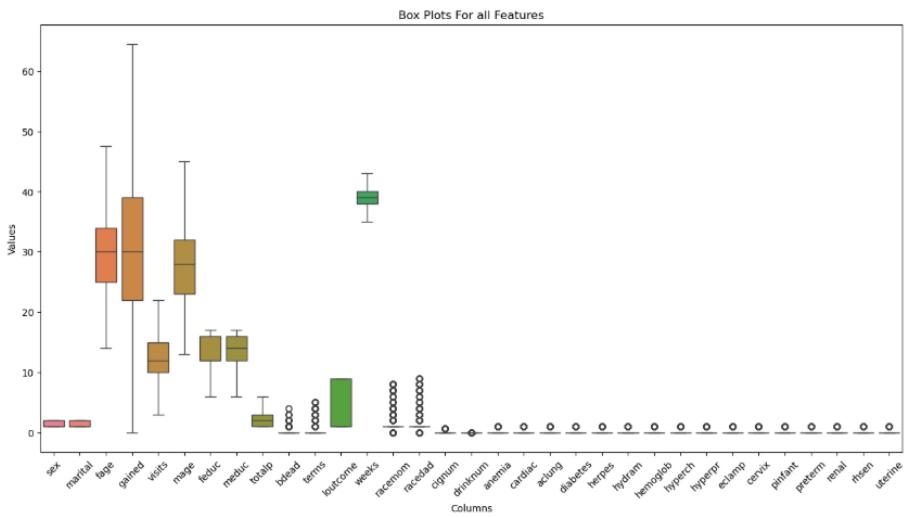
```
[51]: print(len(dataset))
dataset = clip_outliers(dataset, ['fage', 'gained', 'weeks'])
print(len(dataset))

94054
94054
```

```
[52]: # boxplot for all column in the dataset to verify no outlier
plt.figure(figsize=(16, 8))
sns.boxplot(data=dataset.iloc[:, :-1]) # except last/object column
plt.title('Box Plots For all Features')
plt.xticks(rotation=45)
plt.ylabel('Values')
plt.xlabel('Columns')

# plt.savefig('single_boxplot_after_clip_outlier.png')
```

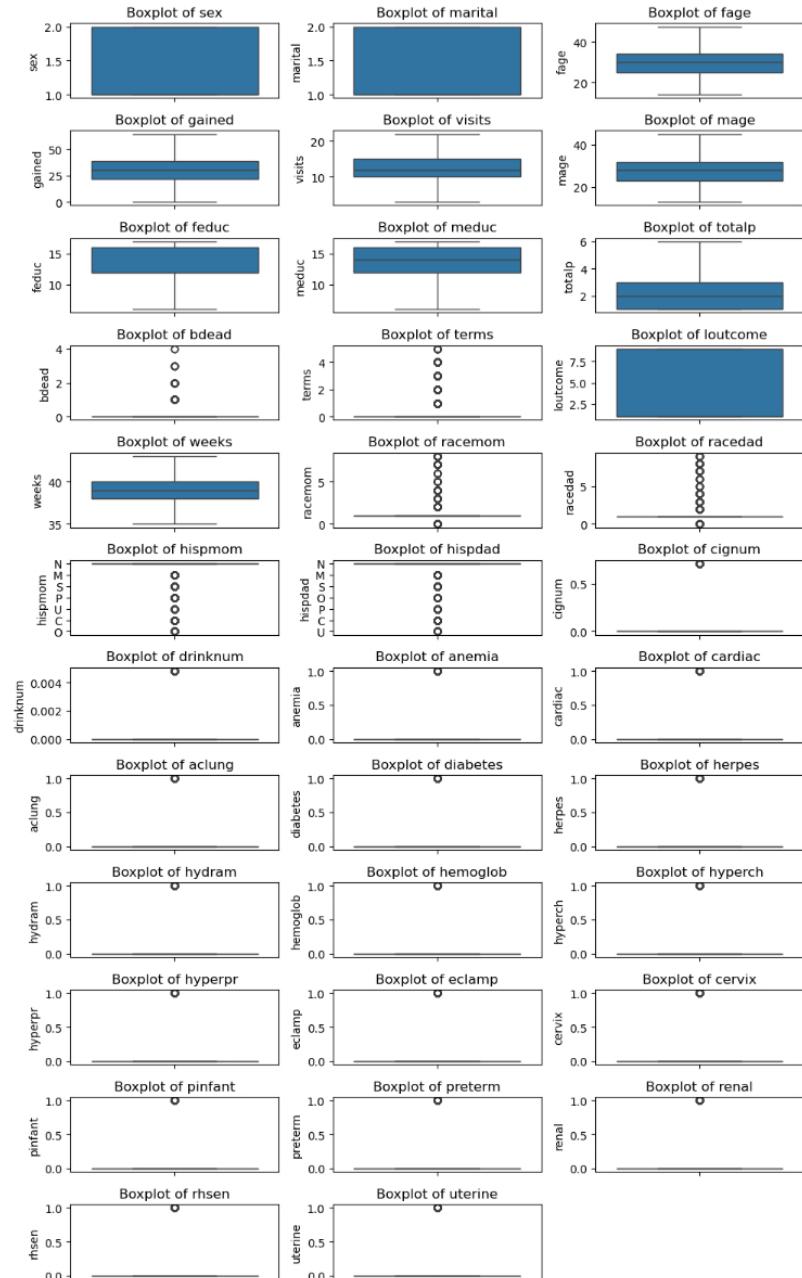
[52]: Text(0.5, 0, 'Columns')



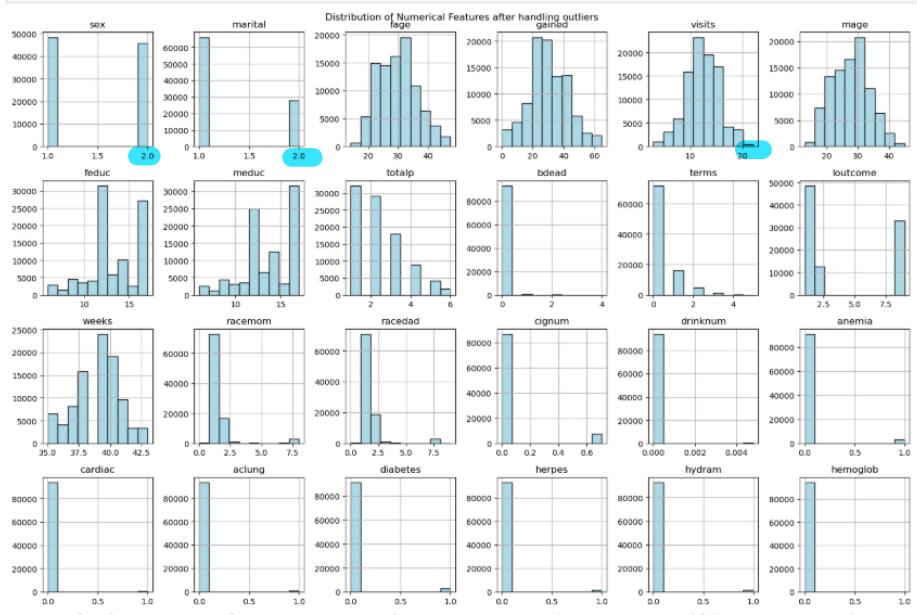
```
[53]: plt.figure(figsize=(10, 16))
```

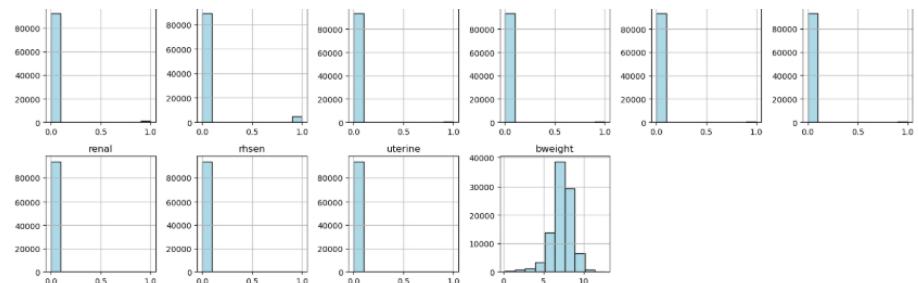
```
# boxplots for each column to verify no outlier
for i, column in enumerate(dataset.columns[:-1], 1): # except the Last/object column
    plt.subplot(12, 3, i) # in 12 rows , 3 columns
    sns.boxplot(data=dataset, y=column)
    plt.title(f'Boxplot of {column}')
plt.tight_layout() # spacing between subplots

# plt.savefig('individual_boxplots_after_outlier_handel.png')
```



```
[54]: # Plot histograms for numerical features
numeric_columns = dataset.select_dtypes(include=['float64', 'int64']).columns
dataset[numeric_columns].hist(bins=10, figsize=(16, 16), color='lightblue', edgecolor='black')
plt.suptitle('Distribution of Numerical Features after handling outliers')
plt.tight_layout()
```





```
[55]: # numeric columns
numeric_columns = dataset.select_dtypes(include=['number']).columns

# calc correlation
correlation_matrix = dataset[numeric_columns].corr()

# correlations with 'BWEIGHT'
correlation_with_bweight = correlation_matrix[['bweight']].sort_values(by='bweight', ascending=False)
correlation_with_bweight.head(15) # most 15 value has influence in bweight
```

[55]: **bweight**

bweight	1.000000
weeks	0.456367
gained	0.178322
visits	0.142505
mage	0.076865
meduc	0.068948
feduc	0.066069
pinfant	0.065061
fage	0.062490
totalp	0.017787
diabetes	0.007971
rhesn	0.002637
herpes	0.001570
cardiac	-0.000638
drinknum	-0.003348

```
[56]: # heatmap
plt.figure(figsize=(6, 10))
sns.heatmap(correlation_with_bweight, annot=True, cmap='Blues', linewidths=0.5, linecolor='black', fmt='.3f')
plt.title('Correlation with bweight')
plt.yticks(rotation=0)

# Save the heatmap
plt.savefig('heatmap_bweight.png')
```

Correlation with bweight

	bweight
bweight	1.000
weeks	0.456
gained	0.178
visits	0.143
mage	0.077
meduc	0.069
feduc	0.066
pinfant	0.065
fage	0.062
totalp	0.018
diabetes	0.008
rhesn	0.003
herpes	0.002
cardiac	-0.001
drinknum	-0.003
aclung	-0.006
renal	-0.006
anemia	-0.010
hemoglobin	-0.012
terms	-0.019
bdead	-0.034
uterine	-0.039
hyperch	-0.042
loutcome	-0.043
hydram	-0.051
cervix	-0.066
preterm	-0.069
eclamp	-0.072
racemom	-0.081
racedad	-0.085
sex	-0.093
cignum	-0.103
marital	-0.111
hyperpr	-0.111
bweight	



We can use this result to drop unnecessary column "not has influence in bweight" or for feature engineering "Create new features or transform existing ones"

Feature Engineering

Creating new features that enhance predictive power can capture more complex relationships in the data and improve model performance.

```
[59]: # new columns can influence in bweight
# average age of the parents might influence the bweight (older parents might have different health conditions or lifestyles)
dataset['avg_parent_age'] = (dataset['mage'] + dataset['fage']) / 2

# total education level of the parents might influence the bweight (more educated parents might follow healthcare)
dataset['total_parent_educ'] = dataset['meduc'] + dataset['feduc']

# provide insight about the mother health during pregnancy
dataset['gained_per_visit'] = dataset['gained'] / dataset['visits']

# to now category: very young or older mothers might have different risks
dataset['mage_category'] = pd.cut(dataset['mage'], bins=[0, 20, 30, 40, 50], labels=['<20', '20-30', '30-40', '>40'])

# more visits during a pregnancy - better health outcomes
dataset['visits_weeks_interaction'] = dataset['visits'] * dataset['weeks']

# combines multiple health-related columns into a single index that represents the overall health status of the mother
# number of health issues
dataset['mom_health_index'] = (
    dataset['anemia'] +
    dataset['cardiac'] +
    dataset['aclung'] +
    dataset['diabetes'] +
    dataset['herpes'] +
    dataset['hydram'] +
    dataset['hemoglobin'] +
```

```

        dataset['hyperch'] +
        dataset['hyperpr'] +
        dataset['eclamp'])

# Categorize the mother's health into groups based on the mom_health_index
dataset['mom_health_category'] = pd.cut(
    dataset['mom_health_index'],
    bins=[-1, 8, 1, 4, float('Inf')]),
    labels=['Good', 'Moderate', 'Poor', 'Very Poor'])

# new features
dataset[['avg_parent_age', 'total_parent_educ', 'gained_per_visit', 'anemia', 'mage_category', 'visits_weeks_interaction', 'cardiac', 'aclung', 'diabetes', 'mom_health_index', 'mom_health_category']]
```

	avg_parent_age	total_parent_educ	gained_per_visit	anemia	mage_category	visits_weeks_interaction	cardiac	aclung	diabetes	mom_health_index	mom_health_category
1	18.5	23.0	4.000000	0	<20	410.0	0	0	0	0.0	Good
2	32.0	32.0	1.142857	0	30-40	546.0	0	0	0	0.0	Good
3	26.5	24.0	2.666667	0	20-30	570.0	0	0	0	0.0	Good
4	20.5	26.0	4.615385	0	<20	520.0	0	0	0	1.0	Moderate
5	21.0	25.0	2.000000	0	20-30	630.0	0	0	0	0.0	Good

identify which features are most correlated with the bweight after make new features

```

[61]: # numeric columns
numeric_columns = dataset.select_dtypes(include=['number']).columns

# calc correlation
correlation_matrix = dataset[numeric_columns].corr()
# correlation_matrix
```

```

[62]: # correlations with bweight
correlation_with_bweight = correlation_matrix[['bweight']].sort_values(by='bweight', ascending=False)
correlation_with_bweight
```

	bweight
bweight	1.000000
weeks	0.456367
visits_weeks_interaction	0.211882
gained	0.178322
visits	0.142505
mage	0.076865
avg_parent_age	0.073778
total_parent_educ	0.072631
meduc	0.068948
feduc	0.066069
pinfant	0.065061
fage	0.062490
gained_per_visit	0.053240
totalp	0.017787
diabetes	0.007971
rhsen	0.002637
herpes	0.001570
cardiac	-0.000638
drinknum	-0.003348
aclung	-0.005654
renal	-0.006348
anemia	-0.010373
hemoglobin	-0.011676
terms	-0.019427
bdead	-0.033505
uterine	-0.038683
hyperch	-0.042369
loutcome	-0.042530
hydram	-0.050683
cervix	-0.065745
preterm	-0.068952
eclamp	-0.071985
racemono	-0.081421
racedad	-0.085485
mom_health_index	-0.092552
sex	-0.092687
cignum	-0.103420
marital	-0.110822
hyperpr	-0.111381

Columns with very low correlation (close to 0) or negative correlation with bweight can be removed Now I will make dataset with features influence on bweight according correlation, and logic knowledge

```

[64]: bweight_dataset = dataset[['
    'weeks', # Strong correlation with bweight
    'visits_weeks_interaction', # Better health outcomes so influence on bweight
    'gained', # Significant for bweight
    'visits', # Indication of prenatal care
    'pinfant', # Historical context, Mother had previous infant 4000+ grams
    'avg_parent_age', # Lifestyle and care on mom and baby health
    'total_parent_educ', # Total education level of parents
    'gained_per_visit', # Weight gained per prenatal visit, insight into pregnancy health
    'mage_category', # Mother's age category
    'mom_health_category', # Risk factors based on age groups
    'cignum', # Risk factors
    'drinknum', # Risk factors
    'diabetes', # Presence of diabetes - dangerous on mom
    'mom_health_index', # Health of mom issues
    'bweight' # Target variable
]]
```

```

bweight_dataset.info()
<class 'pandas.core.frame.DataFrame'
Index: 94854 entries, 1 to 181399
Data columns (total 15 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   weeks            94854 non-null   float64
 1   visits_weeks_interaction 94854 non-null   float64
 2   gained           94854 non-null   float64
 3   visits           94854 non-null   int64  
 4   pinfant          94854 non-null   int64  
 5   avg_parent_age   94854 non-null   float64
 6   mage_category    94854 non-null   float64
 7   mom_health_index 94854 non-null   float64
 8   cignum          94854 non-null   float64
 9   drinknum         94854 non-null   float64
 10  diabetes         94854 non-null   float64
 11  mom_health_cat... 94854 non-null   float64
 12  loutcome         94854 non-null   float64
 13  hydram           94854 non-null   float64
 14  cervix           94854 non-null   float64
 15  preterm          94854 non-null   float64
 16  eclamp           94854 non-null   float64
 17  racemono         94854 non-null   float64
 18  racedad          94854 non-null   float64
 19  outcome           94854 non-null   float64
 20  aclung           94854 non-null   float64
 21  renal             94854 non-null   float64
 22  anemia            94854 non-null   float64
 23  rhsen             94854 non-null   float64
 24  terms             94854 non-null   float64
 25  bdead             94854 non-null   float64
 26  uterine           94854 non-null   float64
 27  herpes            94854 non-null   float64
 28  cardiac            94854 non-null   float64
 29  hyperch            94854 non-null   float64
 30  hyperpr            94854 non-null   float64
 31  sex               94854 non-null   float64
 32  cignum            94854 non-null   float64
 33  drinknum          94854 non-null   float64
 34  diabetes           94854 non-null   float64
 35  mom_health_index  94854 non-null   float64
 36  bweight            94854 non-null   float64
 37  lheight            94854 non-null   float64
 38  racewhite          94854 non-null   float64
 39  raceblack          94854 non-null   float64
 40  raceasian           94854 non-null   float64
 41  raceother           94854 non-null   float64
 42  raceamerindian     94854 non-null   float64
 43  raceasianamerican 94854 non-null   float64
 44  raceblackamerican 94854 non-null   float64
 45  raceamerindianamerican 94854 non-null   float64
 46  raceasianamericanamerican 94854 non-null   float64
 47  raceblackamericanamerican 94854 non-null   float64
 48  raceamerindianamericanamerican 94854 non-null   float64
 49  raceasianamericanamericanamerican 94854 non-null   float64
 50  raceblackamericanamericanamerican 94854 non-null   float64
 51  raceamerindianamericanamericanamerican 94854 non-null   float64
 52  raceasianamericanamericanamericanamerican 94854 non-null   float64
 53  raceblackamericanamericanamericanamerican 94854 non-null   float64
 54  raceamerindianamericanamericanamericanamerican 94854 non-null   float64
 55  raceasianamericanamericanamericanamericanamerican 94854 non-null   float64
 56  raceblackamericanamericanamericanamericanamerican 94854 non-null   float64
 57  raceamerindianamericanamericanamericanamericanamerican 94854 non-null   float64
 58  raceasianamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 59  raceblackamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 60  raceamerindianamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 61  raceasianamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 62  raceblackamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 63  raceamerindianamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 64  raceasianamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 65  raceblackamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 66  raceamerindianamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 67  raceasianamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 68  raceblackamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 69  raceamerindianamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 70  raceasianamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 71  raceblackamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 72  raceamerindianamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 73  raceasianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 74  raceblackamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 75  raceamerindianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 76  raceasianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 77  raceblackamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 78  raceamerindianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 79  raceasianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 80  raceblackamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 81  raceamerindianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 82  raceasianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 83  raceblackamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 84  raceamerindianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 85  raceasianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 86  raceblackamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 87  raceamerindianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 88  raceasianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 89  raceblackamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 90  raceamerindianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 91  raceasianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 92  raceblackamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 93  raceamerindianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 94  raceasianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 95  raceblackamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 96  raceamerindianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 97  raceasianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 98  raceblackamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 99  raceamerindianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 100  raceasianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 101  raceblackamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 102  raceamerindianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 103  raceasianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 104  raceblackamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 105  raceamerindianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 106  raceasianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 107  raceblackamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 108  raceamerindianamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 109  raceasianamericanamericanamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 110  raceblackamericanamericanamericanamericanamericanamericanamerican 94854 non-null   float64
 111  raceamerindianamericanamericanamericanamericanamerican 94854 non-null   float64
 112  raceasianamericanamericanamericanamerican 94854 non-null   float64
 113  raceblackamericanamericanamerican 94854 non-null   float64
 114  raceamerindianamericanamerican 94854 non-null   float64
 115  raceasianamericanamerican 94854 non-null   float64
 116  raceblackamericanamerican 94854 non-null   float64
 117  raceamerindianamerican 94854 non-null   float64
 118  raceasianamerican 94854 non-null   float64
 119  raceblackamerican 94854 non-null   float64
 120  raceamerindian 94854 non-null   float64
 121  raceasian 94854 non-null   float64
 122  raceblack 94854 non-null   float64
 123  raceamerindian 94854 non-null   float64
 124  raceasian 94854 non-null   float64
 125  raceblack 94854 non-null   float64
 126  raceamerindian 94854 non-null   float64
 127  raceasian 94854 non-null   float64
 128  raceblack 94854 non-null   float64
 129  raceamerindian 94854 non-null   float64
 130  raceasian 94854 non-null   float64
 131  raceblack 94854 non-null   float64
 132  raceamerindian 94854 non-null   float64
 133  raceasian 94854 non-null   float64
 134  raceblack 94854 non-null   float64
 135  raceamerindian 94854 non-null   float64
 136  raceasian 94854 non-null   float64
 137  raceblack 94854 non-null   float64
 138  raceamerindian 94854 non-null   float64
 139  raceasian 94854 non-null   float64
 140  raceblack 94854 non-null   float64
 141  raceamerindian 94854 non-null   float64
 142  raceasian 94854 non-null   float64
 143  raceblack 94854 non-null   float64
 144  raceamerindian 94854 non-null   float64
 145  raceasian 94854 non-null   float64
 146  raceblack 94854 non-null   float64
 147  raceamerindian 94854 non-null   float64
 148  raceasian 94854 non-null   float64
 149  raceblack 94854 non-null   float64
 150  raceamerindian 94854 non-null   float64
 151  raceasian 94854 non-null   float64
 152  raceblack 94854 non-null   float64
 153  raceamerindian 94854 non-null   float64
 154  raceasian 94854 non-null   float64
 155  raceblack 94854 non-null   float64
 156  raceamerindian 94854 non-null   float64
 157  raceasian 94854 non-null   float64
 158  raceblack 94854 non-null   float64
 159  raceamerindian 94854 non-null   float64
 160  raceasian 94854 non-null   float64
 161  raceblack 94854 non-null   float64
 162  raceamerindian 94854 non-null   float64
 163  raceasian 94854 non-null   float64
 164  raceblack 94854 non-null   float64
 165  raceamerindian 94854 non-null   float64
 166  raceasian 94854 non-null   float64
 167  raceblack 94854 non-null   float64
 168  raceamerindian 94854 non-null   float64
 169  raceasian 94854 non-null   float64
 170  raceblack 94854 non-null   float64
 171  raceamerindian 94854 non-null   float64
 172  raceasian 94854 non-null   float64
 173  raceblack 94854 non-null   float64
 174  raceamerindian 94854 non-null   float64
 175  raceasian 94854 non-null   float64
 176  raceblack 94854 non-null   float64
 177  raceamerindian 94854 non-null   float64
 178  raceasian 94854 non-null   float64
 179  raceblack 94854 non-null   float64
 180  raceamerindian 94854 non-null   float64
 181  raceasian 94854 non-null   float64
 182  raceblack 94854 non-null   float64
 183  raceamerindian 94854 non-null   float64
 184  raceasian 94854 non-null   float64
 185  raceblack 94854 non-null   float64
 186  raceamerindian 94854 non-null   float64
 187  raceasian 94854 non-null   float64
 188  raceblack 94854 non-null   float64
 189  raceamerindian 94854 non-null   float64
 190  raceasian 94854 non-null   float64
 191  raceblack 94854 non-null   float64
 192  raceamerindian 94854 non-null   float64
 193  raceasian 94854 non-null   float64
 194  raceblack 94854 non-null   float64
 195  raceamerindian 94854 non-null   float64
 196  raceasian 94854 non-null   float64
 197  raceblack 94854 non-null   float64
 198  raceamerindian 94854 non-null   float64
 199  raceasian 94854 non-null   float64
 200  raceblack 94854 non-null   float64
 201  raceamerindian 94854 non-null   float64
 202  raceasian 94854 non-null   float64
 203  raceblack 94854 non-null   float64
 204  raceamerindian 94854 non-null   float64
 205  raceasian 94854 non-null   float64
 206  raceblack 94854 non-null   float64
 207  raceamerindian 94854 non-null   float64
 208  raceasian 94854 non-null   float64
 209  raceblack 94854 non-null   float64
 210  raceamerindian 94854 non-null   float64
 211  raceasian 94854 non-null   float64
 212  raceblack 94854 non-null   float64
 213  raceamerindian 94854 non-null   float64
 214  raceasian 94854 non-null   float64
 215  raceblack 94854 non-null   float64
 216  raceamerindian 94854 non-null   float64
 217  raceasian 94854 non-null   float64
 218  raceblack 94854 non-null   float64
 219  raceamerindian 94854 non-null   float64
 220  raceasian 94854 non-null   float64
 221  raceblack 94854 non-null   float64
 222  raceamerindian 94854 non-null   float64
 223  raceasian 94854 non-null   float64
 224  raceblack 94854 non-null   float64
 225  raceamerindian 94854 non-null   float64
 226  raceasian 94854 non-null   float64
 227  raceblack 94854 non-null   float64
 228  raceamerindian 94854 non-null   float64
 229  raceasian 94854 non-null   float64
 230  raceblack 94854 non-null   float64
 231  raceamerindian 94854 non-null   float64
 232  raceasian 94854 non-null   float64
 233  raceblack 94854 non-null   float64
 234  raceamerindian 94854 non-null   float64
 235  raceasian 94854 non-null   float64
 236  raceblack 94854 non-null   float64
 237  raceamerindian 94854 non-null   float64
 238  raceasian 94854 non-null   float64
 239  raceblack 94854 non-null   float64
 240  raceamerindian 94854 non-null   float64
 241  raceasian 94854 non-null   float64
 242  raceblack 94854 non-null   float64
 243  raceamerindian 94854 non-null   float64
 244  raceasian 94854 non-null   float64
 245  raceblack 94854 non-null   float64
 246  raceamerindian 94854 non-null   float64
 247  raceasian 94854 non-null   float64
 248  raceblack 94854 non-null   float64
 249  raceamerindian 94854 non-null   float64
 250  raceasian 94854 non-null   float64
 251  raceblack 94854 non-null   float64
 252  raceamerindian 94854 non-null   float64
 253  raceasian 94854 non-null   float64
 254  raceblack 94854 non-null   float64
 255  raceamerindian 94854 non-null   float64
 256  raceasian 94854 non-null   float64
 257  raceblack 94854 non-null   float64
 258  raceamerindian 94854 non-null   float64
 259  raceasian 94854 non-null   float64
 260  raceblack 94854 non-null   float64
 261  raceamerindian 94854 non-null   float64
 262  raceasian 94854 non-null   float64
 263  raceblack 94854 non-null   float64
 264  raceamerindian 94854 non-null   float64
 265  raceasian 94854 non-null   float64
 266  raceblack 94854 non-null   float64
 267  raceamerindian 94854 non-null   float64
 268  raceasian 94854 non-null   float64
 269  raceblack 94854 non-null   float64
 270  raceamerindian 94854 non-null   float64
 271  raceasian 94854 non-null   float64
 272  raceblack 94854 non-null   float64
 273  raceamerindian 94854 non-null   float64
 274  raceasian 94854 non-null   float64
 275  raceblack 94854 non-null   float64
 276  raceamerindian 94854 non-null   float64
 277  raceasian 94854 non-null   float64
 278  raceblack 94854 non-null   float64
 279  raceamerindian 94854 non-null   float64
 280  raceasian 94854 non-null   float64
 281  raceblack 94854 non-null   float64
 282  raceamerindian 94854 non-null   float64
 283  raceasian 94854 non-null   float64
 284  raceblack 94854 non-null   float64
 285  raceamerindian 94854 non-null   float64
 286  raceasian 94854 non-null   float64
 287  raceblack 94854 non-null   float64
 288  raceamerindian 94854 non-null   float64
 289  raceasian 94854 non-null   float64
 290  raceblack 94854 non-null   float64
 291  raceamerindian 94854 non-null   float64
 292  raceasian 94854 non-null   float64
 293  raceblack 94854 non-null   float64
 294  raceamerindian 94854 non-null   float64
 295  raceasian 94854 non-null   float64
 296  raceblack 94854 non-null   float64
 297  raceamerindian 94854 non-null   float64
 298  raceasian 94854 non-null   float64
 299  raceblack 94854 non-null   float64
 300  raceamerindian 94854 non-null   float64
 301  raceasian 94854 non-null   float64
 302  raceblack 94854 non-null   float64
 303  raceamerindian 94854 non-null   float64
 304  raceasian 94854 non-null   float64
 305  raceblack 94854 non-null   float64
 306  raceamerindian 94854 non-null   float64
 307  raceasian 94854 non-null   float64
 308  raceblack 94854 non-null   float64
 309  raceamerindian 94854 non-null   float64
 310  raceasian 94854 non-null   float64
 311  raceblack 94854 non-null   float64
 312  raceamerindian 94854 non-null   float64
 313  raceasian 94854 non-null   float64
 314  raceblack 94854 non-null   float64
 315  raceamerindian 94854 non-null   float64
 316  raceasian 94854 non-null   float64
 317  raceblack 94854 non-null   float64
 318  raceamerindian 94854 non-null   float64
 319  raceasian 94854 non-null   float64
 320  raceblack 94854 non-null   float64
 321  raceamerindian 94854 non-null   float64
 322  raceasian 94854 non-null   float64
 323  raceblack 94854 non-null   float64
 324  raceamerindian 94854 non-null   float64
 325  raceasian 94854 non-null   float64
 326  raceblack 94854 non-null   float64
 327  raceamerindian 94854 non-null   float64
 328  raceasian 94854 non-null   float64
 329  raceblack 94854 non-null   float64
 330  raceamerindian 94854 non-null   float64
 331  raceasian 94854 non-null   float64
 332  raceblack 94854 non-null   float64
 333  raceamerindian 94854 non-null   float64
 334  raceasian 94854 non-null   float64
 335  raceblack 94854 non-null   float64
 336  raceamerindian 94854 non-null   float64
 337  raceasian 94854 non-null   float64
 338  raceblack 94854 non-null   float64
 339  raceamerindian 94854 non-null   float64
 340  raceasian 94854 non-null   float64
 341  raceblack 94854 non-null   float64
 342  raceamerindian 94854 non-null   float64
 343  raceasian 94854 non-null   float64
 344  raceblack 94854 non-null   float64
 345  raceamerindian 94854 non-null   float64
 346  raceasian 94854 non-null   float64
 347  raceblack 94854 non-null   float64
 348  raceamerindian 94854 non-null   float64
 349  raceasian 94854 non-null   float64
 350  raceblack 94854 non-null   float64
 351  raceamerindian 94854 non-null   float64
 352  raceasian 94854 non-null   float64
 353  raceblack 94854 non-null   float64
 354  raceamerindian 94854 non-null   float64
 355  raceasian 94854 non-null   float64
 356  raceblack 94854 non-null   float64
 357  raceamerindian 94854 non-null   float64
 358  raceasian 94854 non-null   float64
 359  raceblack 94854 non-null   float64
 360  raceamerindian 94854 non-null   float64
 361  raceasian 94854 non-null   float64
 362  raceblack 94854 non-null   float64
 363  raceamerindian 94854 non-null   float64
 364  raceasian 94854 non-null   float64
 365  raceblack 94854 non-null   float64
 366  raceamerindian 94854 non-null   float64
 367  raceasian 94854 non-null   float64
 368  raceblack 94854 non-null   float64
 369  raceamerindian 94854 non-null   float64
 370  raceasian 94854 non-null   float64
 371  raceblack 94854 non-null   float64
 372  raceamerindian 94854 non-null   float64
 373  raceasian 94854 non-null   float64
 374  raceblack 94854 non-null   float64
 375  raceamerindian 94854 non-null   float64
 376  raceasian 94854 non-null   float64
 377  raceblack 94854 non-null   float64
 378  raceamerindian 94854 non-null   float64
 379  raceasian 94854 non-null   float64
 380  raceblack 94854 non-null   float64
 381  raceamerindian 94854 non-null   float64
 382  raceasian 94854 non-null   float64
 383  raceblack 94854 non-null   float64
 384  raceamerindian 94854 non-null   float64
 385  raceasian 94854 non-null   float64
 386  raceblack 94854 non-null   float64
 387  raceamerindian 94854 non-null   float64
 388  raceasian 94854 non-null   float64
 389  raceblack 94854 non-null   float64
 390  raceamerindian 94854 non-null   float64
 391  raceasian 94854 non-null   float64
 392  raceblack 94854 non-null   float64
 393  raceamerindian 94854 non-null   float64
 394  raceasian 94854 non-null   float64
 395  raceblack 94854 non-null   float64
 396  raceamerindian 94854 non-null   float64
 397  raceasian 94854 non-null   float64
 398  raceblack 94854 non-null   float64
 399  raceamerindian 94854 non
```

```

[65]: bweight_dataset.info()
[65]:
      weeks    visits_weeks_interaction     gained      visits      pinfant avg_parent_age total_parent_educ gained_per_visit      cignum   drinknum
count  94054.000000  94054.000000  94054.000000  94054.000000  94054.000000  94054.000000  94054.000000  94054.000000  94054.000000  94054.000000
mean   38.908393  485.031142  30.026325  12.447317  0.005231  28.767054  26.581304  2.598156  0.056657  0.000011
std    1.874216  130.968335  12.730399  3.259153  0.072137  5.799286  4.952360  1.435349  0.192167  0.000228
min    35.000000  105.000000  0.000000  3.000000  0.000000  13.500000  12.000000  0.000000  0.000000  0.000000
25%   38.000000  400.000000  22.000000  10.000000  0.000000  24.500000  24.000000  1.692308  0.000000  0.000000
50%   39.000000  480.000000  30.000000  12.000000  0.000000  28.500000  26.000000  2.400000  0.000000  0.000000
75%   40.000000  570.000000  39.000000  15.000000  0.000000  33.000000  31.000000  3.250000  0.000000  0.000000
max   43.000000  946.000000  64.500000  22.000000  1.000000  46.250000  34.000000  21.000000  0.708433  0.004827
memory usage: 10.2 MB

[65]: bweight_dataset.describe()

[66]: bweight_dataset['mom_health_category'].value_counts()
[66]: mom_health_category
[66]: Good      79495
[66]: Moderate  12654
[66]: Poor      1985
[66]: Very Poor 0
[66]: Name: count, dtype: int64

[67]: Split data to train and test
[68]: from sklearn.model_selection import train_test_split
[68]: x = bweight_dataset.drop(columns=['bweight'])
[68]: y = bweight_dataset['bweight']
[68]: X_train, X_test, y_train, y_test = train_test_split(x, y, test_size=0.3, random_state=42)

[69]: Data Transformation
[69]: Normalizing or scaling features is important to ensure that all features contribute equally to the model, especially when the features have different scales.

[70]: Standardization and normalization for numeric values
[71]: from sklearn.preprocessing import MinMaxScaler
[71]: # scaler
[71]: scaler = MinMaxScaler()
[71]: # numeric columns
[71]: numeric_columns = X_train.select_dtypes(include=['float64', 'int64']).columns
[71]: # Fit and transform the training data
[71]: X_train[numeric_columns] = scaler.fit_transform(X_train[numeric_columns])
[71]: # transform the testing data
[71]: X_test[numeric_columns] = scaler.transform(X_test[numeric_columns])
[71]: X_train.head(2)

[72]: Encoding Categorical Variables
[73]: from sklearn.preprocessing import OneHotEncoder
[73]: # OneHotEncoder
[73]: encoder = OneHotEncoder(sparse_output=False, handle_unknown='ignore')
[73]: x_train_encoded = encoder.fit_transform(X_train[['mage_category', 'mom_health_category']])
[73]: x_test_encoded = encoder.transform(X_test[['mage_category', 'mom_health_category']])
[74]: encoded_columns = encoder.get_feature_names_out(['mage_category', 'mom_health_category'])
[74]: encoded_columns
[74]: array(['mage_category_>20', 'mage_category_30-40', 'mage_category_<20',
[74]:       'mage_category_>40', 'mom_health_category_Good',
[74]:       'mom_health_category_Moderate', 'mom_health_category_Poor'],
[74]:       dtype=object)
[75]: # x_train_encoded
[76]: # addition - retrieve encoded data as df
[76]: X_train_encoded_df = pd.DataFrame(x_train_encoded, columns=encoded_columns, index=X_train.index)
[76]: X_test_encoded_df = pd.DataFrame(x_test_encoded, columns=encoded_columns, index=X_test.index)
[76]: # drop mage_category , mom_health_category column and add encoded column each represent value
[76]: X_train = X_train.drop(['mage_category', 'mom_health_category'], axis=1).join(X_train_encoded_df)
[76]: X_test = X_test.drop(['mage_category', 'mom_health_category'], axis=1).join(X_test_encoded_df)
[76]: X_train.head(5)

[77]: Train Regressors
[78]: !pip install xgboost
[78]: Requirement already satisfied: xgboost in c:\users\hp\anaconda3\lib\site-packages (2.1.3)
[78]: Requirement already satisfied: numpy in c:\users\hp\anaconda3\lib\site-packages (from xgboost) (1.26.4)
[78]: Requirement already satisfied: scipy in c:\users\hp\anaconda3\lib\site-packages (from xgboost) (1.13.1)
[79]: # dictionary to store RMSE values
[79]: rmse_values = {}

[80]: import numpy as np
[80]: from sklearn.linear_model import Ridge, Lasso
[80]: from scipy.optimize import curve_fit
[80]: from sklearn.linear_model import LinearRegression
[80]: from sklearn.svm import SVR
[80]: from xgboost import XGBRegressor
[80]: from sklearn.ensemble import RandomForestRegressor
[80]: from sklearn.metrics import mean_squared_error

[81]: The Root Mean Squared Error (RMSE) is used to evaluate the accuracy of the models by measuring the difference between the actual and predicted values. It gives a clear indication of how well the model is performing.

[82]: • Linear Regression model "Lasso"
[82]: # Lasso Regression
[82]: lasso_regression = Lasso(alpha=0.01)
[82]: lasso_regression.fit(X_train, y_train)

```

```

y_test_predict_lasso = lasso_regression.predict(X_test)
RMSE_lasso = np.sqrt(mean_squared_error(y_test, y_test_predict_lasso))
rmse_values['Linear Regression - Lasso'] = RMSE_lasso
print(f'Lasso Regression RMSE: {RMSE_lasso}')

```

```

Lasso Regression RMSE: 1.1417017145014268

```

```

[83]: results_df = pd.DataFrame({
    'Actual Value': y_test,
    'Predicted Value': y_test_predict_lasso
})

```

```

results_df.head(10)

```

	Actual Value	Predicted Value
35304	9.1250	7.745116
77984	9.4375	6.370752
61265	8.2500	7.355441
49828	7.4375	7.189246
61851	8.3750	7.147922
6765	7.6250	7.097600
16206	6.0625	6.244758
64608	7.5000	8.189014
87712	6.9375	6.945761
33846	5.5000	6.658820

```

[84]: plt.figure(figsize=(10, 6))

```

```

# actual values
plt.scatter(y_test, y_test_predict_lasso, color='dodgerblue', alpha=0.6, edgecolor='k', s=100, label='Predicted Values')

```

```

# Link actual and predicted values with lines
for actual, predicted in zip(y_test, y_test_predict_lasso):
    plt.plot([actual, actual], [actual, predicted], 'k--', alpha=0.3)

```

```

# plot perfect fit line
plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='red', linewidth=2, label='Perfect Fit Line')

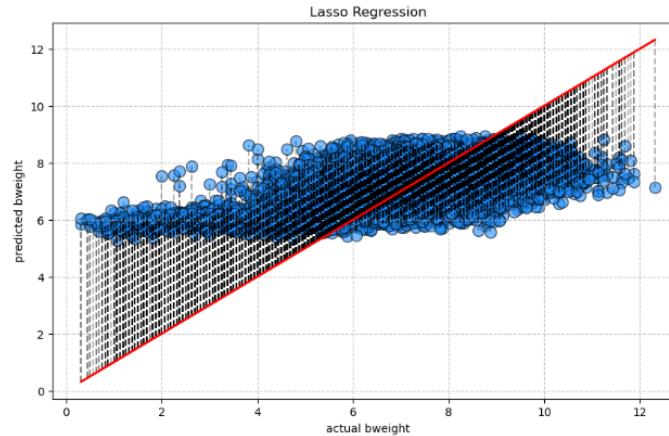
plt.grid(True, linestyle='--', alpha=0.7)
plt.xlabel('actual bweight')
plt.ylabel('predicted bweight')
plt.title('Lasso Regression')

```

```

[84]: Text(0.5, 1.0, 'Lasso Regression')

```



• Linear Regression model "Ridge"

```

[86]: # Ridge Regression
ridge_regression = Ridge()
ridge_regression.fit(X_train, y_train)
y_test_predict_ridge = ridge_regression.predict(X_test)
RMSE_ridge = np.sqrt(mean_squared_error(y_test, y_test_predict_ridge))
rmse_values['Linear Regression - Ridge'] = RMSE_ridge
print(f'Ridge Regression RMSE: {RMSE_ridge}')

```

```

Ridge Regression RMSE: 1.1260093264996474

```

```

[87]: results_df = pd.DataFrame({
    'Actual Value': y_test,
    'Predicted Value': y_test_predict_ridge
})

```

```

results_df.head(10)

```

	Actual Value	Predicted Value
35304	9.1250	7.752874
77984	9.4375	6.743556
61265	8.2500	7.343751
49828	7.4375	7.155635
61851	8.3750	7.110418
6765	7.6250	6.925508
16206	6.0625	5.873396
64608	7.5000	8.263143
87712	6.9375	6.971915
33846	5.5000	6.805177

```

[88]: plt.figure(figsize=(10, 6))

```

```

# actual values
plt.scatter(y_test, y_test_predict_ridge, color='dodgerblue', alpha=0.6, edgecolor='k', s=100, label='Predicted Values')

```

```

# Link actual and predicted values with lines
for actual, predicted in zip(y_test, y_test_predict_ridge):
    plt.plot([actual, actual], [actual, predicted], 'k--', alpha=0.3)

```

```

# fit Line
plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='red', linewidth=2, label='Perfect Fit Line')

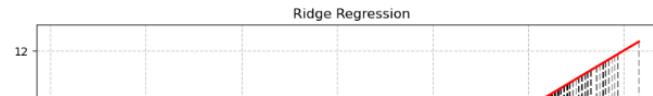
plt.xlabel('actual bweight')
plt.ylabel('predicted bweight')
plt.title('Ridge Regression')

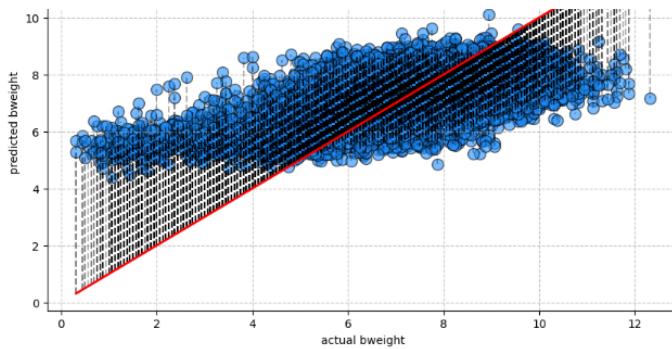
```

```

plt.grid(True, linestyle='--', alpha=0.7)

```





• `fit_curve` for linear regression

```
[90]: # train the Linear regression model
model = LinearRegression()
model.fit(X_train, y_train)

# equation
print("Linear Regression Equation:")
intercept = model.intercept_
coefficients = model.coef_
features = X_train.columns

equation_parts = [(f"(coef) * ({feature})" for coef, feature in zip(coefficients, features))]
equation = f"{intercept} + " + " + ".join(equation_parts)

equation
```

Linear Regression Equation:

$$2333847119327.875 + 3.8123461829624214 \cdot \text{weeks} + -12.110715741227978 \cdot \text{visits_weeks_interaction} + 0.9158277380170403 \cdot \text{gained} + 11.06629692438277 \cdot \text{visits}$$

$$+ 2333847119327.875 + 3.8123461829624214 \cdot \text{weeks} + -12.110715741227978 \cdot \text{visits_weeks_interaction} + 0.9158277380170403 \cdot \text{gained} + 11.06629692438277 \cdot \text{visits}$$

$$+ 0.37222515285 * \text{pinfant} + 0.37322515285 * \text{avg_parent_age} + -0.060398331722995226 * \text{total_parent_educ} + 0.12695086388385395 * \text{gained_per_visit}$$

$$+ -0.3972266481529195 * \text{cignum} + -0.012927862185556211 * \text{drinknum} + 0.4758520895164058 * \text{diabetes} + -0.392888172408453 * \text{mom_health_index} + -241891473896.8$$

$$-2091955645426.1292 * \text{mom_health_category_Good} + -2091955645426.2793 * \text{mom_health_category_Moderate} + -2091955645426.4304 * \text{mom_health_category_Poor}$$

```
[91]: # predict values using test set
y_test_predict_linear = model.predict(X_test)

# RMSE
RMSE_linear = np.sqrt(mean_squared_error(y_test, y_test_predict_linear))
rmse_values['fit_curve - Linear Regression'] = RMSE_linear
print(f"Linear Regression RMSE: {RMSE_linear}")

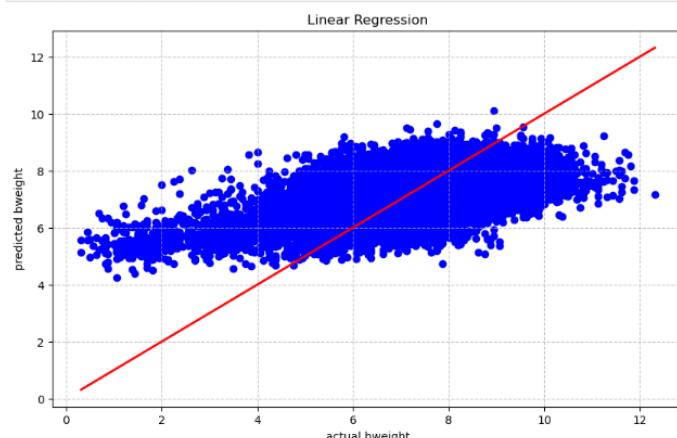
Linear Regression RMSE: 1.1251808101200589
```

```
[92]: results_df = pd.DataFrame({
    'Actual Value': y_test,
    'Predicted Value': y_test_predict_linear
})

results_df.head(10)
```

	Actual Value	Predicted Value
35304	9.1250	7.739746
77984	9.4375	6.742676
61265	8.2500	7.349121
49828	7.4375	7.150391
61851	8.3750	7.115723
6765	7.6250	6.929199
16206	6.0625	5.825195
64608	7.5000	8.266602
87712	6.9375	6.982422
33846	5.5000	6.802734

```
[93]: # Scatter plot of actual vs predicted values
plt.figure(figsize=(10, 6))
plt.scatter(y_test, y_test_predict_linear, color='blue', label='Actual vs Predicted')
plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='red', linewidth=2, label='Perfect Fit Line')
plt.xlabel('actual bweight')
plt.ylabel('predicted bweight')
plt.title('Linear Regression')
plt.grid(True, linestyle='--', alpha=0.7)
```



```
[94]: plt.figure(figsize=(10, 6))

# scatterm actual values
plt.scatter(y_test, y_test_predict_linear, color='dodgerblue', alpha=0.6, edgecolor='k', s=100, label='Predicted Values')

# Link actual and predicted values with lines
for actual, predicted in zip(y_test, y_test_predict_linear):
    plt.plot([actual, actual], [actual, predicted], 'k--', alpha=0.3)

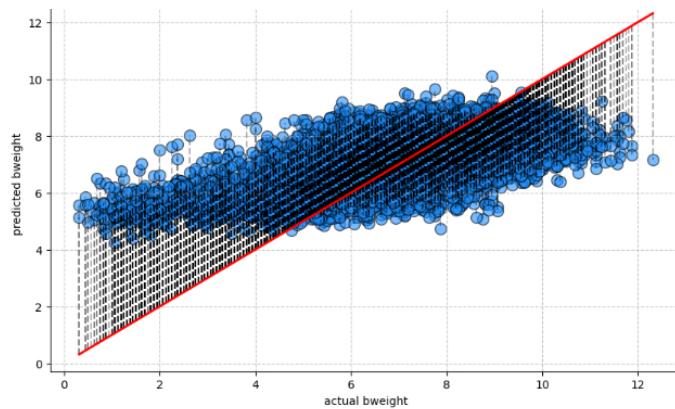
# perfect fit line
plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='red', linewidth=2, label='Perfect Fit Line')

plt.grid(True, linestyle='--', alpha=0.7)
plt.xlabel('actual bweight')
plt.ylabel('predicted bweight')
plt.title('Linear Regression')

plt.figure(figsize=(10, 6))
```

```
[94]: Text(0.5, 1.0, 'Linear Regression')
```

Linear Regression



• `fit_curve` for non-linear regression

using Polynomial model

$$\text{Polynomial model: } (y = a_0 + a_1x + \dots + a_nx^n)$$

```
[97]: from sklearn.preprocessing import PolynomialFeatures
from sklearn.pipeline import make_pipeline

# polynomial regression for different degrees
def evaluate_polynomial_regression(X_train, X_test, y_train, y_test, max_degree):
    rmse_values = []
    for degree in range(1, max_degree + 1):
        # Polynomial regression model
        poly = PolynomialFeatures(degree=degree)
        poly_regression = make_pipeline(poly, LinearRegression())
        poly_regression.fit(X_train, y_train)

        # Predict values
        y_test_pred_poly = poly_regression.predict(X_test)

        # RMSE
        RMSE_poly = np.sqrt(mean_squared_error(y_test, y_test_pred_poly))
        rmse_values.append((degree, RMSE_poly))

    # Print results dataframe
    results_df = pd.DataFrame({
        'Actual Value': y_test,
        'Predicted Value': y_test_pred_poly
    })

    print(results_df.head(10))
    print(f"Polynomial Regression (Degree {degree}) RMSE: {RMSE_poly}")

    # Plot result
    plt.figure(figsize=(10, 6))
    plt.scatter(y_test, y_test_pred_poly, color='dodgerblue', alpha=0.6, edgecolor='k', s=100, label='Predicted Values')
    plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='red', linewidth=2, label='Perfect Fit Line')
    plt.xlabel('actual Bweight')
    plt.ylabel('predicted Bweight')
    plt.title('Nonlinear - Polynomial Regression (Degree {degree})')
    plt.grid(True, linestyle='--', alpha=0.7)
    plt.show()

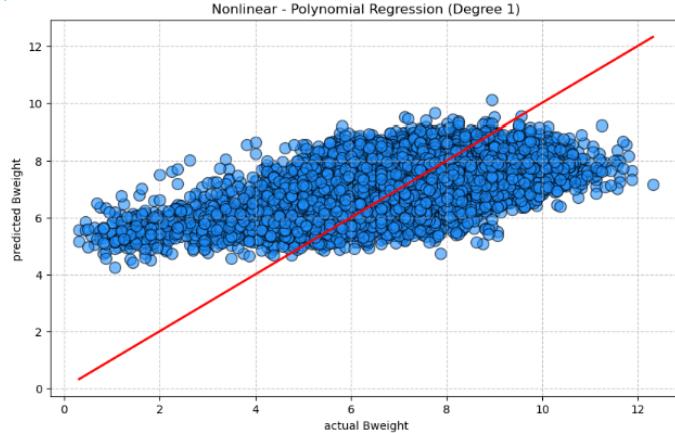
    return rmse_values
```

```
[98]: # Polynomial regression for degrees from 1 to 3
max_degree = 3
rmse_values_poly = evaluate_polynomial_regression(X_train, X_test, y_train, y_test, max_degree)

# Find the degree with the minimum RMSE
best_degree, best_rmse = min(rmse_values_poly, key=lambda x: x[1])
rmse_values['Fit Curve - Nonlinear Regression - polynomial Model'] = best_rmse
print(f"The best degree for polynomial regression is: {best_degree} with RMSE: {best_rmse}")
```

	Actual Value	Predicted Value
35304	9.1250	7.740479
77984	9.4375	6.742676
61265	8.2500	7.349365
49828	7.4375	7.150391
61851	8.3750	7.115723
6765	7.6250	6.930908
16206	6.0625	5.826566
64608	7.5000	8.266846
87712	6.9375	6.982910
33846	5.5000	6.882979

n = 1 same to linear

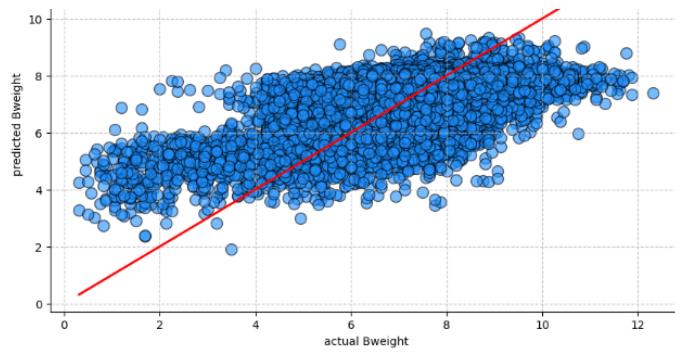


	Actual Value	Predicted Value
35304	9.1250	7.876118
77984	9.4375	6.322567
61265	8.2500	7.737377
49828	7.4375	7.383186
61851	8.3750	7.383308
6765	7.6250	7.258278
16206	6.0625	6.438278
64608	7.5000	7.748340
87712	6.9375	7.228153
33846	5.5000	6.712306

Polynomial Regression (Degree 2) RMSE: 1.0615413585607085

Nonlinear - Polynomial Regression (Degree 2)



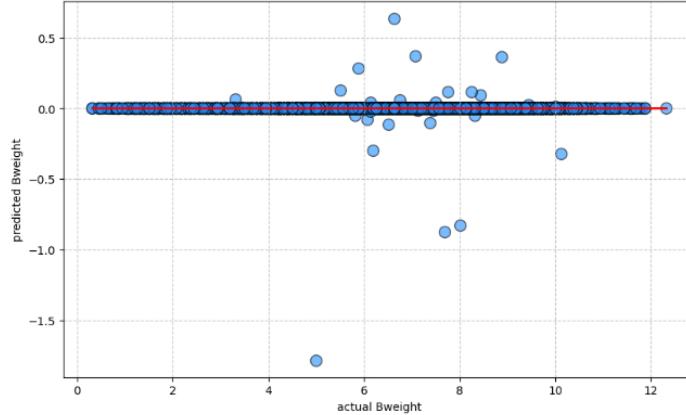


Actual Value Predicted Value

35304	9.1250	7.879822
77984	9.4375	6.513367
61265	8.2500	7.601467
49828	7.4375	7.337494
61851	8.3750	7.400879
6765	7.6250	7.328186
16206	6.0625	6.543152
64608	7.5000	7.596924
87712	6.9375	7.244293
33846	5.5000	6.545776

Polynomial Regression (Degree 3) RMSE: 141864725.78008565

Nonlinear - Polynomial Regression (Degree 3)



The best degree for polynomial regression is: 2 with RMSE: 1.0615413585607085

using Power model

Power model:

$$(y = ax^b)$$

```
[100]: def power_func(x, a, b):
    return a * np.power(x, b)

[101]: coefficients = []
equations = []
for feature in X_train.columns:
    # find the best parameters for the power function for each feature
    params, _ = curve_fit(power_func, X_train[feature], y_train)
    coefficients.append(params)
    a, b = params
    equation = f'y = ({a:.4f}) * ({feature})^{({b:.4f})}'
    equations.append(equation)

# Print the equations for each feature
print("Power Regression Model Equations:")
for eq in equations:
    print(eq)
```

```
C:\Users\HP\AppData\Local\Temp\ipykernel_14296\2598982826.py:5: OptimizeWarning: Covariance of the parameters could not be estimated
  params, _ = curve_fit(power_func, X_train[feature], y_train)
Power Regression Model Equations:
y = 1.7418 * weeks^0.0000
y = 1.8979 * visits_weeks_interaction^0.0000
y = 7.7790 * gained^0.0000
y = 1.8016 * visits^0.0000
y = 8.4681 * pifant^1.0000
y = 2.6187 * avg_parent_age^0.0000
y = 1.4723 * total_parent_educ^0.0000
y = 4.6523 * gained_per_visit^0.0000
y = 6.8242 * cignum^1.0000
y = 7.2865 * drinknum^1.0000
y = 7.3354 * diabetes^1.0000
y = 2.9642 * mom_health_index^0.0000
y = 7.2775 * mage_category_29-38^1.0000
y = 7.3678 * mage_category_39-48^1.0000
y = 7.8251 * mage_category_>28^1.0000
y = 7.1746 * mage_category_>49^1.0000
y = 7.3224 * mom_health_category_Good^1.0000
y = 7.6225 * mom_health_category_Moderate^1.0000
y = 6.7507 * mom_health_category_Poor^1.0000
```

```
[102]: # Predict values
y_test_predict_power = np.zeros_like(y_test)
# predicted target values using the fitted power functions for each feature
for i, feature in enumerate(X_test.columns):
    a, b = coefficients[i]
    y_test_predict_power += power_func(X_test[feature], a, b) / len(X_test.columns)

# RMSE
RMSE_power = np.sqrt(mean_squared_error(y_test, y_test_predict_power))
rmse_values['Fit_Curve - Nonlinear Regression - Power Model'] = RMSE_power
print(f"Power Regression RMSE: {RMSE_power}")

Power Regression RMSE: 5.00085928382347
```

```
[103]: results_df = pd.DataFrame({
    'Actual Value': y_test,
    'Predicted Value': y_test_predict_power
})

results_df.head(10)
```

Actual Value Predicted Value

35304	9.1250	1.871269
77984	9.4375	2.397571
61265	8.2500	1.871269
49828	7.4375	1.866517
61851	8.3750	1.866517
6765	7.6250	2.011499
16206	6.0625	2.212402

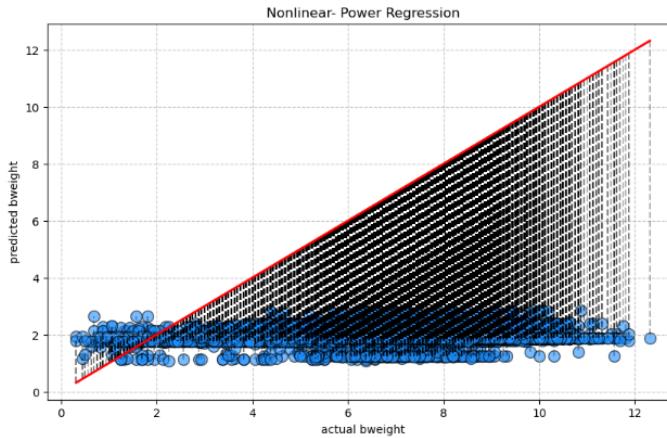
```
[104]: plt.figure(figsize=(10, 6))

# Scatter plot of actual vs predicted values
plt.scatter(y_test, y_test_predict_power, color='dodgerblue', alpha=0.6, edgecolor='k', s=100, label='Predicted Values')

# Connect actual and predicted values with lines
for actual, predicted in zip(y_test, y_test_predict_power):
    plt.plot([actual, actual], [actual, predicted], 'k--', alpha=0.3)

# fit line
plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='red', linewidth=2, label='Perfect Fit Line')

plt.xlabel('actual bweight')
plt.ylabel('predicted bweight')
plt.title('Nonlinear Power Regression')
plt.grid(True, linestyle='--', alpha=0.7)
```



• Support Vector Machine (SVM) Regressor

```
[106]: # SVM Regressor
svm_regression = SVR(kernel='rbf')
# fit train data
svm_regression.fit(X_train, y_train)
# predict test data values
y_test_predict_svm = svm_regression.predict(X_test)
RMSE_svm = np.sqrt(mean_squared_error(y_test, y_test_predict_svm))
rmse_values['Nonlinear Regression - SVM'] = RMSE_svm
print("SVM Regressor RMSE: ", RMSE_svm)
```

Uses the rbf kernel to effectively model non-linear relationships while minimizing error within a margin of tolerance.

```
[107]: results_df = pd.DataFrame({
    'Actual Value': y_test,
    'Predicted Value': y_test_predict_svm
})

results_df.head(10)
```

	Actual Value	Predicted Value
35304	9.1250	7.866741
77984	9.4375	6.416038
61265	8.2500	7.705368
49828	7.4375	7.323704
61851	8.3750	7.372983
6765	7.6250	7.329343
16206	6.0625	6.206452
64608	7.5000	7.563322
87712	6.9375	7.207808
33846	5.5000	6.782067

```
[108]: plt.figure(figsize=(10, 6))

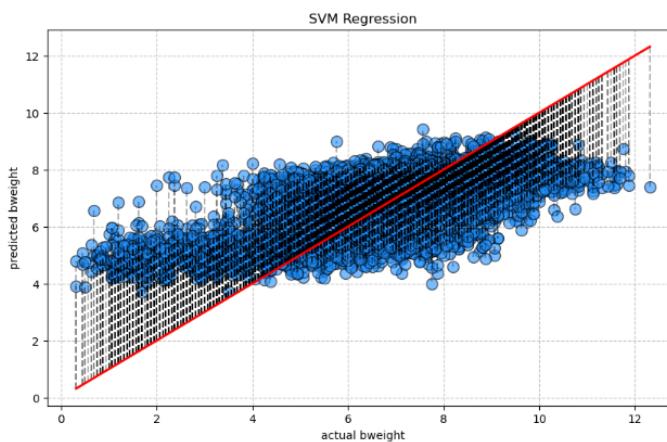
# Scatter plot of actual vs predicted values
plt.scatter(y_test, y_test_predict_svm, color='dodgerblue', alpha=0.6, edgecolor='k', s=100, label='Predicted Values')

# Connect actual and predicted values with lines
for actual, predicted in zip(y_test, y_test_predict_svm):
    plt.plot([actual, actual], [actual, predicted], 'k--', alpha=0.3)

# perfect fit line
plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='red', linewidth=2, label='Perfect Fit Line')

plt.xlabel('actual bweight')
plt.ylabel('predicted bweight')
plt.title('SVM Regression')
plt.grid(True, linestyle='--', alpha=0.7)
```

```
[108]: Text(0.5, 1.0, 'SVM Regression')
```



XGBoost Regressor

Known for high performance and efficiency, it uses gradient boosting to optimize model performance through additive models.

```
[110]: # XGBoost Regressor

# Convert feature names to strings and replace any special characters to make XGBRegressor can deal with
X_train.columns = [str(col).replace('[', '_').replace(']', '_').replace('<', '_') for col in X_train.columns]
X_test.columns = X_train.columns

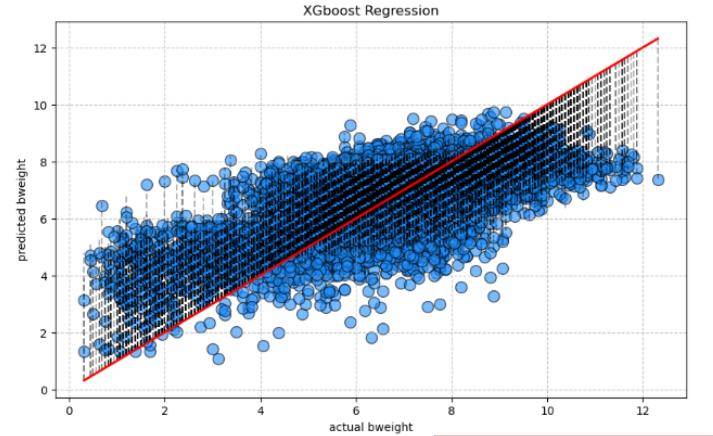
xboost_regression = XGBRegressor()
# fit train data values
xboost_regression.fit(X_train, y_train)
# predict test data values
y_test_predict_xgboost = xboost_regression.predict(X_test)
RMSE_xgboost = np.sqrt(mean_squared_error(y_test, y_test_predict_xgboost))
rmse_values['XGBoost Regressor'] = RMSE_xgboost
print("XGBoost Regressor RMSE: " + str(RMSE_xgboost))

XGBoost Regressor RMSE: 1.050762440712795

[111]: results_df = pd.DataFrame({
    'Actual Value': y_test,
    'Predicted Value': y_test_predict_xgboost
})

results_df.head(10)

[111]:   Actual Value Predicted Value
0  35304.0  9.1250
1  77984.0  9.4375
2  61265.0  8.2500
3  49828.0  7.4375
4  61851.0  8.3750
5  6765.0   7.6250
6  16206.0  6.0625
7  64608.0  7.5000
8  87712.0  6.9375
9  33846.0  5.5000
```



Random Forest Regressor (with 50 estimators)

Combines multiple decision trees to handle non-linear relationships and interactions, providing robust predictions.

```
[114]: rf_regression = RandomForestRegressor(n_estimators=50, random_state=42)
# fit train data values
rf_regression.fit(X_train, y_train)
# predict test data values
y_test_predict_rf = rf_regression.predict(X_test)
RMSE_rf = np.sqrt(mean_squared_error(y_test, y_test_predict_rf))
rmse_values['Random Forest Regressor'] = RMSE_rf
print("Random Forest Regressor RMSE: " + str(RMSE_rf))

Random Forest Regressor RMSE: 1.0967006944764879

[115]: results_df = pd.DataFrame({
    'Actual Value': y_test,
    'Predicted Value': y_test_predict_rf
})

results_df.head(10)

[115]:   Actual Value Predicted Value
0  35304.0  8.027917
1  77984.0  6.996562
2  61265.0  8.102500
3  49828.0  7.573750
4  61851.0  6.692500
5  6765.0   7.250000
6  16206.0  5.597500
7  64608.0  7.655000
8  87712.0  7.265937
9  33846.0  6.733125
```

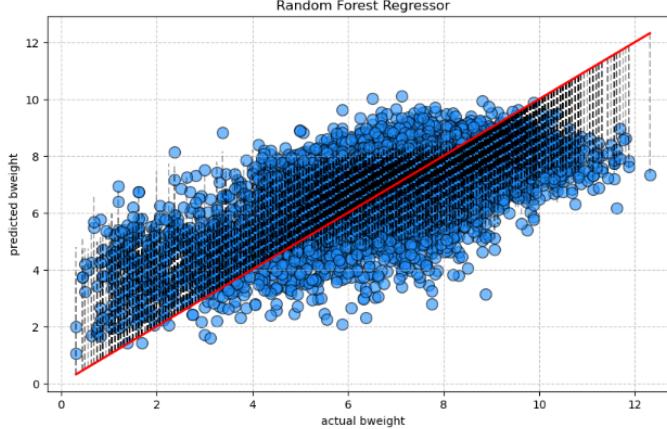


```

plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='red', linewidth=2, label='Perfect Fit Line')
plt.xlabel('actual bwheight')
plt.ylabel('predicted bwheight')
plt.grid(True, linestyle='--', alpha=0.7)
plt.title('Random Forest Regressor')

```

[116]: Text(0.5, 1.0, 'Random Forest Regressor')



Evaluate Regressors

```

[118]: # Create a mapping dictionary for new names
key_mapping = {
    'Linear Regression - Lasso': 'Lasso',
    'Linear Regression - Ridge': 'Ridge',
    'fit_curve - Nonlinear Regression - polynomial Model': 'Polynomial Regression',
    'fit_curve - Nonlinear Regression - Power Model': 'Power Regression',
    'Nonlinear Regression - SVM': 'SVM',
    'XGBoost Regressor': 'XGBoost',
    'Random Forest Regressor': 'Random Forest'
}

rmse_values_plot = {key_mapping.get(key, key): value for key, value in rmse_values.items()}

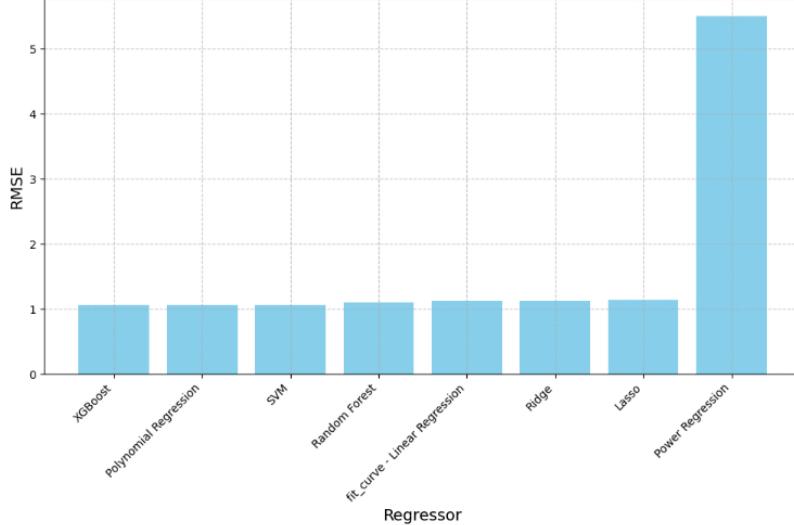
rmse_df = pd.DataFrame(list(rmse_values_plot.items()), columns=['Regressor', 'RMSE'])

rmse_df = rmse_df.sort_values(by='RMSE').reset_index(drop=True)

plt.figure(figsize=(12, 6))
plt.bar(rmse_df['Regressor'], rmse_df['RMSE'], color='skyblue')
plt.xlabel('Regressor', fontsize=14)
plt.ylabel('RMSE', fontsize=14)
plt.title('Comparison of Regressors using RMSE', fontsize=16)
plt.xticks(rotation=45, ha='right')
plt.grid(True, linestyle='--', alpha=0.7)

```

Comparison of Regressors using RMSE



```

[119]: rmse_df = pd.DataFrame(list(rmse_values.items()), columns=['Regressor', 'RMSE'])
rmse_df_sorted = rmse_df.sort_values(by='RMSE')
rmse_df_sorted

```

	Regressor	RMSE
6	XGBoost Regressor	1.059749
3	fit_curve - NonLinear Regression - polynomial ...	1.061541
5	Nonlinear Regression - SVM	1.063992
7	Random Forest Regressor	1.096701
2	fit_curve - Linear Regression	1.125181
1	Linear Regression - Ridge	1.126009
0	Linear Regression - Lasso	1.141702
4	fit_curve - Nonlinear Regression - Power Model	5.500986

```

[120]: fig, axs = plt.subplots(3, 2, figsize=(15, 12), sharex=True, sharey=True)

# Linear Regression (Lasso)
axs[0, 0].scatter(y_test, y_test_predict_lasso, color='blue', alpha=0.6, edgecolor='k', s=50, label='Lasso')
axs[0, 0].plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='black', linewidth=2, linestyle='--')
axs[0, 0].set_title('Lasso')

# Linear Regression (Ridge)
axs[0, 1].scatter(y_test, y_test_predict_ridge, color='green', alpha=0.6, edgecolor='k', s=50, label='Ridge')
axs[0, 1].plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='black', linewidth=2, linestyle='--')
axs[0, 1].set_title('Ridge')

# NonLinear - Power Regression
axs[1, 0].scatter(y_test, y_test_predict_power, color='orange', alpha=0.6, edgecolor='k', s=50, label='Power Regression')
axs[1, 0].plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='black', linewidth=2, linestyle='--')
axs[1, 0].set_title('Power Regression')

# NonLinear SVM Regressor
axs[1, 1].scatter(y_test, y_test_predict_svm, color='red', alpha=0.6, edgecolor='k', s=50, label='SVM')

```

```

axs[1, 1].plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='black', linewidth=2, linestyle='--')
axs[1, 1].set_title('SVM')

# XGBoost Regressor
axs[2, 0].scatter(y_test, y_test_predict_xgboost, color='purple', alpha=0.6, edgecolor='k', s=50, label='XGBoost')
axs[2, 0].plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='black', linewidth=2, linestyle='--')
axs[2, 0].set_title('XGBoost')

# Random Forest Regressor
axs[2, 1].scatter(y_test, y_test_predict_rf, color='cyan', alpha=0.6, edgecolor='k', s=50, label='Random Forest')
axs[2, 1].plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='black', linewidth=2, linestyle='--')
axs[2, 1].set_title('Random Forest')

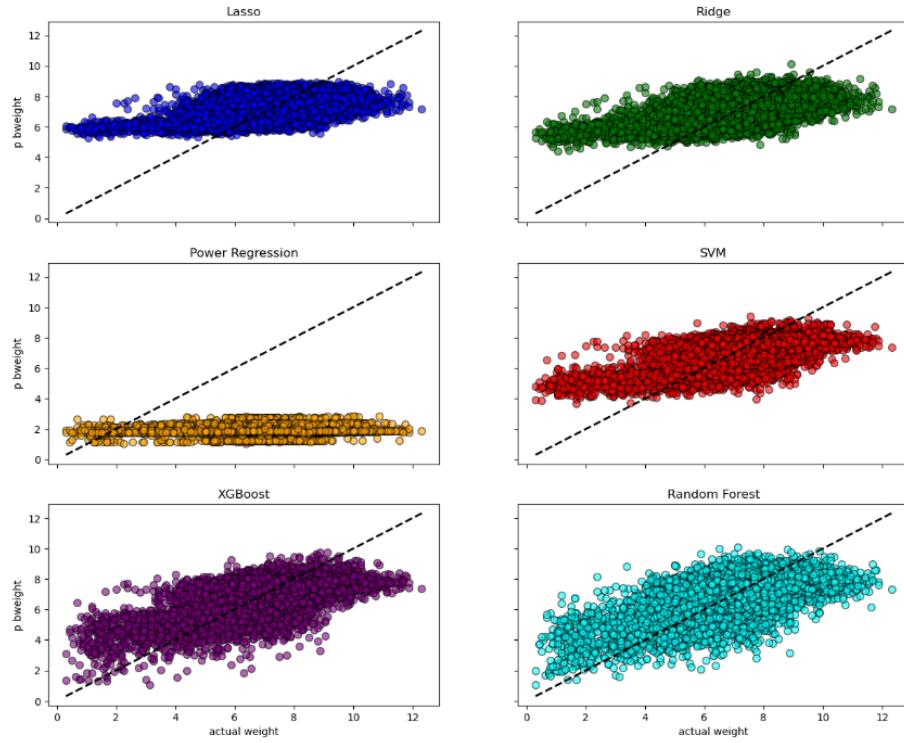
# common labels
for ax in axs.flat:
    ax.set(xlabel='actual weight', ylabel='p bweight')
    ax.label_outer()

fig.suptitle('Comparison of Different Regressors', fontsize=16)

# Show legend
handles, labels = axs[0, 0].get_legend_handles_labels()
plt.savefig('all_Regressors.png')

```

Comparison of Different Regressors



[]: