# Validation using Cross-Validation method

Name: Likhil Naik Vislavath

## The Issues

1) Use the validation set method as described on pages 198-200 of the text to split the data into two random halves, using one half as the training set and the remaining half as the test set.

2) Use leave-one-outcross-validation (LOOCV), as described on pages 200-202 of the text, to test the linear model.

3) Use k-fold cross-validation, with k = 10, as described on pages 203-206 of the text, to test the linear model.

## Findings

The multivariate linear regression model fitted to data from the Babiesweight.xls file was tested using three different cross-validation methods. The obtained results are as follows: R-squared of 0.276 for the validation set method, mean R-squared of 0.283 for the LOOCV method, and mean R-squared of 0.280 for the 10-fold cross-validation method. These results indicate that the model has moderate predictive ability for birthweight based on the five predictor variables: Gestation, Age, Height, Weight, and Smoke. However, the R-squared values are relatively low, suggesting that the model can only explain a small portion of the variation in birthweight. The consistency of the R-squared values across the three cross-validation methods indicates that the model is not overfitting to the training data and can generalize well to new data. Nevertheless, the results imply that factors beyond the five predictor variables likely influence birthweight, limiting the model's usefulness for predicting birthweight to some extent.

### Discussion

### Appendix A : Method

The dataset used for predicting the birthweight of babies included five predictor variables: Gestation, Age, Height, Weight, and Smoke. The data was initially in .xls format and was imported into Jupyter notebook. In order to enhance the model's predictive accuracy,

multiple predictor variables were utilized and issues in the data pre-processing and model fitting stages were resolved. To evaluate the model's predictive accuracy, various methods such as the validation set, LOOCV, and k-fold cross-validation were employed.

## Appendix B: Results

Despite identifying height and smoking as the most strongly correlated predictor variables for birth weight, the new model only achieved an R-squared value of 0.031. Moreover, the model's predictive accuracy was found to be inadequate as indicated by the high MSE on the validation set, and the high mean absolute error on both LOOCV and k-fold cross-validation methods. Therefore, it is not advisable to use this model for practical purposes, and further enhancements may be necessary.

```
                         OLS Regression Results
==============================================================================
Dep. Variable:            Birthweight   R-squared:                       0.031
Model:                            OLS   Adj. R-squared:                  0.027
Method:                 Least Squares   F-statistic:                     7.754
Date:                Fri, 31 Mar 2023   Prob (F-statistic):           3.41e-07
Time:                        15:32:36   Log-Likelihood:                -5322.8
No. Observations:                1236   AIC:                         1.066e+04
Df Residuals:                    1230   BIC:                         1.069e+04
Df Model:                           5
Covariance Type:            nonrobust
==============================================================================
                 coef    std err          t      P>|t|      [0.025      0.975]
------------------------------------------------------------------------------
const          81.8104      7.947     10.294      0.000      66.219      97.402
Gestation       0.0128      0.007      1.874      0.061      -0.001       0.026
Age             0.0704      0.079      0.886      0.376      -0.086       0.226
Height          0.5256      0.122      4.311      0.000       0.286       0.765
Weight         -0.0058      0.004     -1.345      0.179      -0.014       0.003
Smoke          -1.9890      0.562     -3.542      0.000      -3.091      -0.887
==============================================================================
Omnibus:                       13.075   Durbin-Watson:                   2.048
Prob(Omnibus):                  0.001   Jarque-Bera (JB):               17.582
Skew:                          -0.118   Prob(JB):                     0.000152
Kurtosis:                       3.534   Cond. No.                     5.40e+03
==============================================================================

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 5.4e+03. This might indicate that there are
strong multicollinearity or other numerical problems.
```

**figure 1: OLS regression results**

## Appendix C: Code

```
import pandas as pd
```

```python
import numpy as np

from sklearn.linear_model import LinearRegression

from sklearn.model_selection import train_test_split, cross_val_score,
LeaveOneOut, KFold

# Load data from Excel file

data = pd.read_excel("Babies_weight.xls")

X = data[['Gestation', 'Age', 'Height', 'Weight', 'Smoke']]

y = data['Birthweight']

# Fit multivariate linear regression model

model = LinearRegression().fit(X, y)

# Use validation set method to split data into training and test sets

X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.5,

random_state=42)

model.fit(X_train, y_train)

y_pred = model.predict(X_test)

print("Validation set method R-squared:", model.score(X_test, y_test))

# Use LOOCV to test linear model

loocv = LeaveOneOut()

scores = cross_val_score(model, X, y, cv=loocv)

print("LOOCV mean R-squared:", np.mean(scores))

# Use k-fold cross-validation to test linear model

kfold = KFold(n_splits=10, shuffle=True, random_state=42)

scores = cross_val_score(model, X, y, cv=kfold)

print("10-fold cross-validation mean R-squared:", np.mean(scores))
```

**References:**

https://chat.openai.com/chat