

5	1989-08	8.09	20.5	27.4	34.3	0.28	1	76.10
6	1989-09	111.88	16.4	23.3	30.2	28.30	1	0.26
7	1989-10	0.00	10.2	17.9	25.7	0.40	1	28.30
8	1989-11	0.00	4.8	13.4	22.1	0.02	1	0.40
9	1989-12	15.88	-1.2	6.8	14.8	1.23	1	0.02

note that you may not all the values for each case, depending on your best model.

Your "Model Building: Part 1" section must include your answers for the following questions:

- ```

\
str(round(slope,2)) + 'x' + \
) + \

```

```
plt.xlabel('Precipitation(in.)')
plt.ylabel('Flowrate(cfs)')
plt.legend(['Monthly Precipitation', 'Model Prediction'])
```

- ```

Method: Least Squares      F-statistic: 225.2
Date, 03 May 2022         Prob (F-statistic): 4.56e-49
Time: 23:07:27             Log-Likelihood: -1912.2
No. Observations: 373      AIC: 3828.
DF Residuals: 371         BIC: 3856.
DF Model: 1
Covariance Type: nonrobust

=====
               coef      std err          t      P>|t|    [0.025    0.975]
-----+-----
Intercept    -18.3377      2.989      -6.153    0.000    -24.193    -12.474
PPT           0.7412      0.049      15.006    0.000      0.644      0.838
=====
Omnibus: 566.188   Durbin-Watson:      1.989
Prob(Omnibus): 0.000   Jarque-Bera (JB): 145657.838
Skew: 7.667         Prob(JB):      0.00
Kurtosis: 96.320     Cond. No.      85.0
=====

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```
-
- ```

In [57]: # statspackage model - cubic

df['x2'] = df['PPT']**2 # add a column of x^2
df['x3'] = df['PPT']**3 # add a column of x^3

Initialise and fit linear regression model using 'statsmodels'
cubemodel = smf.ols('flow_cfs ~ PPT+X2+X3', data=df) # model object constructor syntax
cubemodel = cubemodel.fit()

Predict values
cubeppt_pred = cubemodel.predict()
print(cubemodel.summary())

Intercept = cubemodel.params[0]
slope1 = cubemodel.params[1]
slope2 = cubemodel.params[2]

```

```

rsquare = cubemodel.rsquared
RMSE = math.sqrt(cubemodel.mse_total)

Plot
titleline = 'Precipitation effect on flow rate \n' + \
 'y = ' + str(round(intercept,2)) + ' + ' + str(round(slope1,3)) + 'x' + ' + ' + str(round(slope2,3)) + 'x^2' + ' + ' + str(round(slope3,3)) + 'x^3' + \
 '\n R squared = ' + str(round(Rsquare,3)) + \
 '\n RMSE = ' + \
 str(round(RMSE,2))

plt.figure(figsize=(12, 6))
plt.plot(df['PPT'], df['Flow_cfs'], 'o') # scatter plot showing actual data
plt.plot(df['PPT'], cubicppt_pred, markers='s', color='r', linewidth=0) # regression line
plt.xlabel('Precipitation(in.)')
plt.ylabel('Flowrate(cfs)')
plt.legend(['Monthly Precipitation', 'Model Prediction'])
plt.title(titleline)
plt.show();

=====
OLS Regression Results
=====
Dep. Variable: Flow_cfs R-squared: 0.872
Model: OLS Adj. R-squared: 0.871
Method: Least Squares F-statistic: 834.7
Date: Tue, 03 May 2022 Prob (F-statistic): 5.12e-164
Time: 23:07:27 Log-Likelihood: -1617.9
No. Observations: 373 AIC: 3244
Df Residuals: 369 BIC: 3259
Df Model: 3
Covariance Type: nonrobust
=====
coef std err t P>|t| [0.025 0.975]

Intercept -5.8596 1.936 -3.026 0.003 -9.667 -2.052
PPT 0.8368 0.108 7.749 0.000 0.624 1.049
XX -0.0147 0.001 -10.548 0.000 -0.017 -0.012
XXX 8.409e-05 4.51e-06 18.643 0.000 7.52e-05 9.3e-05

Omnibus: 191.211 Durbin-Watson: 1.918
Prob(Omnibus): 0.000 Jarque-Bera (JB): 4658.266
Skew: 1.606 Prob(JB): 0.00
Kurtosis: 20.012 Cond. No. 3.03e+08
=====

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 3.03e+08. This might indicate that there are
strong multicollinearity or other numerical problems.

Precipitation effect on flow rate
y = -5.86 + 0.837x + -0.015x^2 + 0.0001x^3
R squared = 0.872
RMSE = 51.73


```

```
from sklearn.linear_model import LinearRegression

instantiate the model (using the default parameters)
linreg = LinearRegression()

fit the model with data
linreg.fit(x_train,y_train)

run predictor on x_test
y_pred=linreg.predict(x_test)

In [60]: # use statspackage model cubic for testing

df[['XX']]=df[['PPT']]**2 # add a column of X^2
df[['XXX']]=df[['PPT']]**3 # add a column of X^2

Initialise and fit linear regression model using 'statsmodels'
cubemodel = smf.ols('Flow.cfs ~ PPT+XX+XXX', data=df) # model object constructor syntax
cubemodel = cubemodel.fit()
Predict values
cubicppt_pred = cubemodel.predict()

intercept = cubemodel.params[0]
slope1 = cubemodel.params[1]
slope2 = cubemodel.params[2]
slope3 = cubemodel.params[3]
Rsquare = cubemodel.rsquared
RMSE = math.sqrt(cubemodel.mse_total)

Apply equation from cubic model as function for interactive input
def cubic(i, intercept, slope1, slope2, slope3):
 ans = intercept + (slope1 * i) + (slope2 * math.pow(i, 2)) + (slope3 * math.pow(i,3))
 if(ans<0):
 ans = 0.0
 return(ans)
```

```

input data interface
i = input("Enter requested precipitation to test: ")
i = np.float64(i)
flowrate = cubic(i, intercept, slope1, slope2, slope3)
print("Expected flowrate is: ", "%-2f" % format(flowrate), " cfs")

Enter requested precipitation to test: 11
Expected flowrate is: 1.48 cfs

```

- What are the most important assumptions in your modeling?
  - The most important assumption are to neglect the extreme values in the dataset that lead to outliers because the goal of the model is to create a pattern from the trends seen in the data, and by that eliminating the outliers, the data is more consistent with a wider range of data.
- Is it beneficial to use the streamflow recorded in the previous step (a lagged streamflow value) as an input feature? why?
  - No because based on the correlation values, using the lagged streamflow value had a very minuscule correlation to the flow\_cfs, being the lowest out of all of the parameters.
- Which parameter (between precipitation and temperature values) would be a better predictor for streamflow at the station of study? why?
  - Using the precipitation parameter was a better predictor for streamflow due to how the precipitation has a bigger r-squared value compared to temperature, due to the r-squared value of precipitation being 0.378 out of 1, which is better than the r-squared value of temperature being 0.015.
- Is there a specific range of streamflow values that are harder to capture accurately for your data models? If yes, what range and why?
  - The range of precipitation values from 0.0 to 8.0 were the hardest to capture accurately for the data models, likely due to the majority of those streamflow cases having a flowrate of 0 and the outliers had flowrates.

## Model Building: Part 2 | Forecasting flow states

In this part, the goal is to make models to predict whether the Colorado River's flowstate is in the "Flow" or the "No-Flow" state. Then, evaluate their performance using appropriate goodness-of-fit measures, and analyze the outcomes. Use the first 75% of the dataset for training your models and the remaining 25% for testing.

- Add a column to the dataframe for the flow state: It should be 0 when the flowrate is equal to 0 and 1 when the flowrate is non-zero.
- Build 3 data models that you see appropriate for predicting the flow state. (Please note that these models should be unique and this uniqueness can be defined based on using different algorithms, inputs, or both.)
- Assess data model quality (decide which model is best)
- Build the input data interface for using the "best" model
- Using your best model determine the estimated flow state for the hydro-meteorological conditions in the table below: [ppt[mm][mean][min] | :---|---|---| : 0.0113 0.99 0.015 0.015 0.95 0.95 0.75 0.1 2.2] 0.10 0.00 0.1 0.1 0.80 0.06 0.40 0.0

[0.0180 0.00 0.40 0.0] data that you may not all of the flow rates for each case, depending on your best model.

- ```
model = model.fit()
print(model.summary())
```

```
str(round(RMSE, 2))
```

```
plt.figure(figsize=(12, 6))
plt.plot(df['PPT'], df['Flow?'], 'o') # scatter plot showing actual data
plt.plot(df['PPT'], y_pred, marker='s', color='r', linewidth=8) # regression line
plt.xlabel('Precipitation(in.)')
```

- | | | | |
|-------------------|-----|------|-------|
| No. observations: | 373 | AIC: | 403.4 |
| Df Residuals: | 370 | BIC: | 415.2 |
| Df Model: | 2 | | |

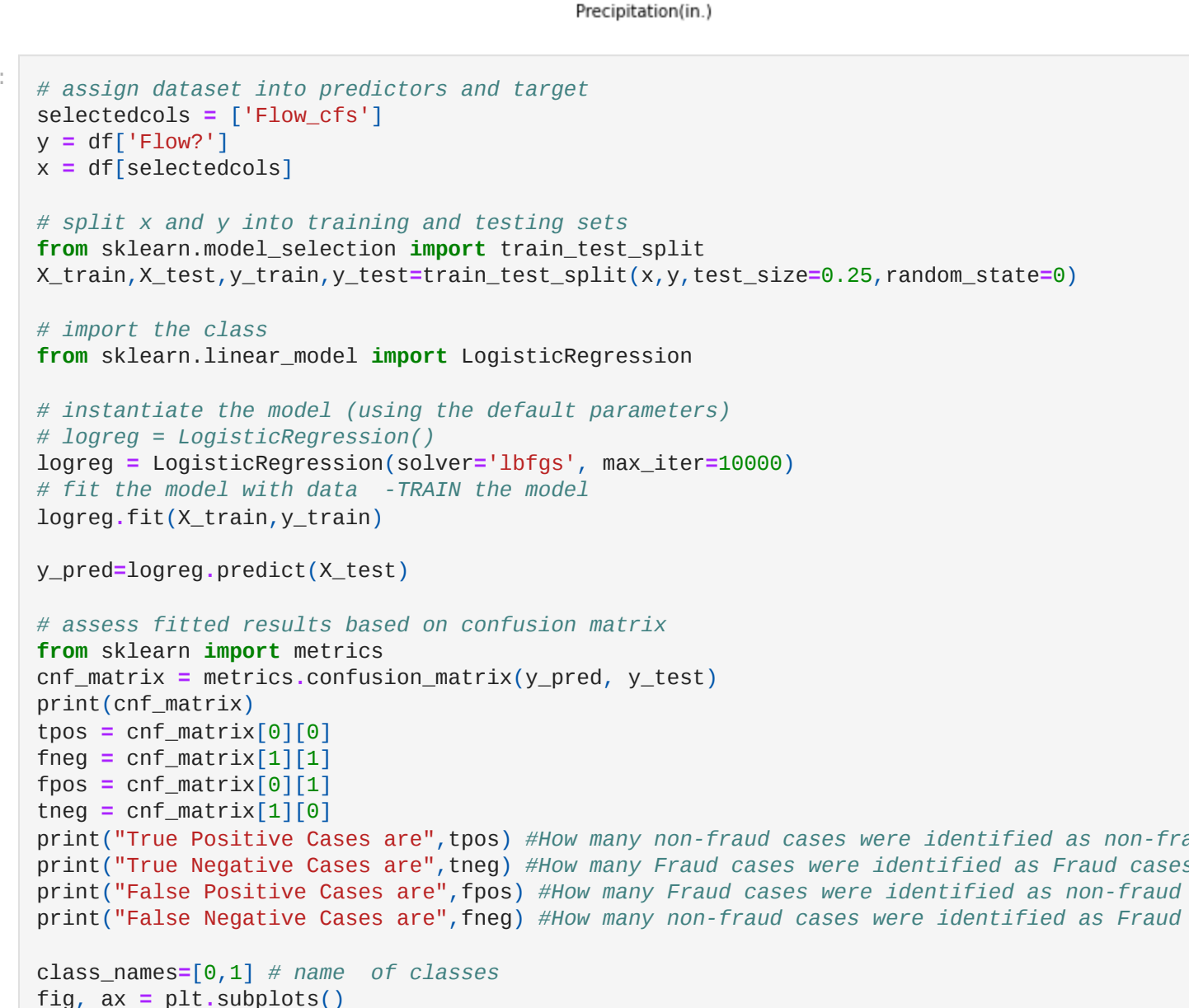
```
=====
            coef      std err      t      Pr>|t|      [0.025      0.975]
-----+-----
Intercept      0.4416      0.052     13.935      0.000      0.379      0.584
Flow_CFs       0.0015      0.001     -2.332      0.005      -0.003      +0.009
PPT            0.0060      0.001      9.513      0.000      0.005      0.007
-----+-----
Omnibus:                104.088      Durbin-Watson:                1.293
Prob(Omnibus):          0.000      Jarque-Bera (JB):          35.731
Skew:                   -0.336      Prob(JB):                  1.74e-88
Kurtosis:               1.641      Cond. No.                  107.
=====
```

Notes:

[3] Standard Errors assume that the covariance matrix of the errors is correctly specified.

Precipitation effect on flow rate

R squared = 0.221
RMSE = 51.73



```
plt.xticks(tick_marks, class_names)
plt.yticks(tick_marks, class_names)

# create heatmap
import seaborn as sns
sns.heatmap(pd.DataFrame(conf_matrix), annot=True, cmap='YlGnBu', ftext='g')
ax.xaxis.set_label_position('top')
plt.tight_layout()
plt.title('Confusion matrix', y=1.1)
plt.ylabel('Predicted label')
plt.xlabel('Actual label')

print("Accuracy:", metrics.accuracy_score(y_test, y_pred))
print("Precision:", metrics.precision_score(y_test, y_pred))
print("Recall:", metrics.recall_score(y_test, y_pred))
print("F1-score:", metrics.f1_score(y_test, y_pred))
from sklearn.metrics import classification_report

print(classification_report(y_test, y_pred))

[[34  7]
 [ 0 53]]
True Positive Cases are 34
True Negative Cases are 0
False Positive Cases are 7
False Negative Cases are 53
Accuracy: 0.905531914890617
Precision: 1.0
```

	precision	recall	f1-score	support
0	0.83	1.00	0.91	34
1	1.00	0.88	0.94	60
accuracy			0.93	94
macro avg	0.91	0.94	0.92	94
weighted avg	0.94	0.93	0.93	94

Confusion matrix

```

In [64]: # sklearn - K-Nearest Neighbors model

# split data into predictors and target
selectedcols = ['Flow_cfs']
X = df[selectedcols]
y = df['Flow?']

# split data into training and testing with 25:75 split
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25)

# import KNeighborsClassifier Class
from sklearn.neighbors import KNeighborsClassifier
classifier = KNeighborsClassifier(n_neighbors=4)
classifier.fit(X_train, y_train)

# make predictions on test data

```

```

y_pred = classifier.predict(X_test)

# assess fitted results using confusion matrix
from sklearn.metrics import classification_report, confusion_matrix
print(confusion_matrix(y_test, y_pred))
print(classification_report(y_test, y_pred))

# apply to confusion matrix
cm = pd.DataFrame(confusion_matrix(y_test, y_pred))
sns.heatmap(cm, annot=True)
plt.title('Confusion matrix', y=1.1)
plt.xlabel('Predicted label')
plt.ylabel('Actual label')

[[31 0]
 [ 2 61]]

      precision    recall  f1-score   support

      0       0.94       1.00       0.97        31
      1       1.00       0.97       0.98        63

 accuracy          0.97          0.98          0.98          94
 macro avg          0.97          0.98          0.98          94
 weighted avg          0.98          0.98          0.98          94

Text(0.5, 15.0, 'Actual label')

```

```

In [65]: # finding the best possible K-value

error = []

# Calculating error for K values between 1 and 50
# In each iteration the mean error for predicted values of test set is calculated and
# the result is appended to the error list.
for i in range(1, 50):
    knn = KNeighborsClassifier(n_neighbors=i)
    knn.fit(X_train, y_train)
    pred_i = knn.predict(X_test)
    error.append(np.mean(pred_i != y_test))

# plot the error values against K values:
plt.figure(figsize=(12, 6))
plt.plot(range(1, 50), error, color='red', linestyle='dashed', markers='o',
        markerfacecolor='blue', markersize=10)
plt.title('Error Rate K Value')
plt.xlabel('K value')

```

```
plt.ylabel('Mean Error')
plt.minorticks_on()
plt.grid(which='both')
```

K Value	Mean Error
0	0.000
1	0.000
2	0.000
3	0.021
4	0.021
5	0.021
6	0.021
7	0.021
8	0.021
9	0.021
10	0.021
11	0.021
12	0.021
13	0.021
14	0.021
15	0.021
16	0.021
17	0.021
18	0.021
19	0.021
20	0.021
21	0.021
22	0.021
23	0.021
24	0.021
25	0.021
26	0.021
27	0.021
28	0.021
29	0.021
30	0.021
31	0.021
32	0.032
33	0.032
34	0.032
35	0.032
36	0.032
37	0.032
38	0.032
39	0.032
40	0.032
41	0.032
42	0.032
43	0.032
44	0.032
45	0.032
46	0.032
47	0.032
48	0.032
49	0.032
50	0.032

```
# Input data interface
ppt = input("Enter requested precipitation to test: ")
ppt = np.float64(ppt)
flowrate = cubic(ppt, intercept, slope1, slope2, slope3)
print("Expected flowrate is: {}".format(flowrate), " cfs")
print(classifier.predict([flowrate]))
```

Enter requested precipitation to test: 200
Expected flowrate is: 246.72 cfs
[1]

- What can be some applications of this kind of streamflow state modeling?

- The streamflow state modelling can be used for agriculture and public safety in order to determine how much water would affect the area of study.