

Importing libraries

import numpy as np
import pandas as pd
import statistics
import math
from matplotlib import pyplot as plt
import statsmodels.formula.api as smf

Reading and defining a dataframe of the csv file

df = pd.read_csv("ColoradoRiverData.csv")
df.head()

DatePPTTminTavTmaxdeltSMIETFlow_cfsFlow?

01988-0316.092.812.321.819.0197.90731.4298452.681
11988-0417.637.917.026.118.2299.71061.4270572.331
21988-0563.441.14.21.428.714.61367.616104.9317534.561
31988-0658.0118.625.833.014.41496.658149.37451718.201
41988-0791.9820.526.632.712.22446.668161.75150676.101

Using describe function to obtain important statistical information and basic information

df.describe()

PPTTminTavTmaxdeltSMIETFlow_cfsFlow?

count373.000000373.000000373.000000373.000000373.000000373.000000373.000000373.000000373.000000
mean42.48321710.45655817.77104625.08766814.633099862.2586680.64093030.1535060.678284
std48.9488248.1826697.8965887.71573561.868985104.12544363.16559951.7293980.467762
min0.000000-4.6000003.3000008.8000009.9000000.0000001.9347900.0000000.000000
25%0.9100002.50000010.50000018.20000013.400000126.10000019.6479380.0000000.000000
50%30.70000010.40000018.10000025.90000014.000000428.28000065.1742000.0000001.000000
75%62.85000018.40000025.30000031.80000016.0000001389.024000138.8432357.6000001.000000
max251.24000024.30000031.50000038.80000019.4000005672.650000218.125353709.9000001.000000

Finding the amount of columns and rows in dataframe

df.shape

(373, 10)

Using the correlation function to get the correlation

corr = df.corr(method = 'pearson')
corr

PPTTminTavTmaxdeltSMIETFlow_cfsFlow?

PPT1.0000000.3776170.3175930.249051-0.6181470.9399770.2984870.0145880.452687
Tmin0.3776171.0000000.9937090.973441-0.3555010.5772040.5761490.1512650.234497
Tav0.3175930.9937091.0000000.992938-0.2466630.5200830.9751570.1234330.197389
Tmax0.2490510.9734410.9929381.000000-0.1320760.4688360.9628770.0923450.179355
delt-0.618147-0.355501-0.246663-0.1320761.000000-0.565551-0.266853-0.277924-0.382443
SMI0.9399770.5772040.5200830.468836-0.5655511.0000000.5022350.5759030.423881
ET0.2984870.9741840.9751570.962877-0.2668530.5022351.0000000.1012660.184756
Flow_cfs0.0145880.1512650.1234330.092345-0.2779240.5759030.1012661.0000000.173555
Flow?0.4526870.2344970.1973890.156013-0.3824430.4238810.1847560.1735551.000000

Determining how many rows of data have a flow value of 0 and move them into its own dataframe

NoFlowPeriods = df[df["Flow?"] < 0.5]
NoFlowPeriods

DatePPTTminTavTmaxdeltSMIETFlow_cfsFlow?

191989-106.5011.019.427.716.7126.10031.4298790.00
202018-110.834.012.621.317.310.45828.0284360.00
211989-126.00-4.63.311.215.819.8002.1513300.00
271990-0611.2121.929.336.714.8328.453180.00
331990-1218.24-2.05.813.615.6105.7826.3277310.00
...
3592018-028.06-0.18.517.017.168.51012.8788510.00
3602018-0313.156.514.823.116.6194.62044.7786750.00
3612018-040.696.916.426.019.111.31657.3463630.00
3702019-016.00-0.16.813.813.940.8008.7402910.00
3712019-023.711.69.417.215.634.87415.6133470.00
120 rows x 10 columns

Using the describe function to view basic and statistical information of the new dataframe of rows only with non-flowperiods

NoFlowPeriods.describe()

PPTTminTavTmaxdeltSMIETFlow_cfsFlow?

count120.000000120.000000120.000000120.000000120.000000120.000000120.000000120.000000120.000000
mean12.3045837.67416715.51083323.35333315.757957238.8430250.5677090.000.00
std17.2256528.2395478.0416757.9489861.459497344.11249264.494110.000.00
min0.000000-4.6000003.30000010.70000012.800000270.64800028.8857250.4800000.00
25%2.3800000.5750008.57500016.65000014.47500022.84960013.4314160.000.00
50%8.6900005.35000013.35000021.90000015.650000133.06600031.7259710.000.00
75%20.67000013.22500021.02500028.87500016.4000001078.095000147.84105013.9000001.00
max104.37000024.30000031.50000038.80000019.4000002035.210000218.1253530.000.00

Determining how many rows of data have a flow value of 1 and move them into its own dataframe

FlowingPeriods = df[df["Flow?"] > 0.5]
FlowingPeriods

DatePPTTminTavTmaxdeltSMIETFlow_cfsFlow?

01988-0316.092.812.321.819.0197.90731.4298452.681
11988-0417.637.917.026.118.2299.71061.4270572.331
21988-0563.441.14.21.428.714.61367.616104.9317534.561
31988-0658.0118.625.833.014.41496.658149.37451718.201
41988-0791.9820.526.632.712.22446.668161.75150676.101
...
3662018-0977.5717.723.228.711.01799.624106.11028022.201
3672018-10169.5111.516.922.410.92864.71954.78820839.201
3682018-118.101.89.116.314.573.71015.0391820.101
3692018-1230.350.16.613.012.9200.3108.1022000.011
3722019-0219.364.310.917.513.2211.02424.9428570.151
253 rows x 10 columns

Using the describe function to view basic and statistical information of the new dataframe of rows only with non-flowperiods

FlowingPeriods.describe()

PPTTminTavTmaxdeltSMIETFlow_cfsFlow?

count253.000000253.000000253.000000253.000000253.000000253.000000253.000000253.000000253.000000
mean46.92011911.77628518.84303325.91027714.1339921157.84988185.51686919.3822331.0
std44.9546667.8491787.6106537.4787511.6389961090.78423161.05209531.8765990.0
min0.000000-2.4000003.5000008.8000009.9000000.0000002.4530930.0100001.0
25%22.0200004.50000012.00000019.60000012.800000270.64800028.8857250.4800001.0
50%45.18000012.90000020.20000027.50000014.000000838.10100090.3594872.7600001.0
75%75.25000019.40000025.70000032.20000015.4000001498.260000136.95970410.7750001.0
max251.24000022.90000030.10000037.20000019.0000005672.650000218.125353709.9000001.0

Highlighting rows that were recorded only in the 1990s and moving it into its own dataframe

The1990s = df.iloc[22:142]
The1990s

DatePPTTminTavTmaxdeltSMIETFlow_cfsFlow?

221990-0116.490.48.318.215.8136.86712.7979402.941
231990-0235.392.410.618.916.5374.49819.6479380.541
241990-0347.185.912.318.912.9580.06831.42964510.001
251990-0497.9810.117.424.714.61703.11264.23180251.501
261990-0512.5314.421.628.914.5270.648106.8158622.321
...
1271999-0845.7722.229.629.914.81254.792188.7402411.531
1281999-0942.6616.923.930.914.01019.574112.2358671.201
1291999-109.169.918.026.215.2164.80061.8018590.000
1401999-110.006.014.723.417.40.00037.6426920.000
1411999-120.00-0.37.918.016.30.00011.4283270.000
120 rows x 10 columns

Using the describe function to obtain basic and statistical information of the 1990s dataframe

The1990s.describe()

PPTTminTavTmaxdeltSMIETFlow_cfsFlow?

count120.000000120.000000120.000000120.000000120.000000120.000000120.000000120.000000120.000000
mean40.18425010.22866717.61333325.00666714.780000803.83693378.86228511.1919170.616667
std40.3964088.0514397.7672657.5816531.718282955.28689962.67510932.7869150.488237
min0.000000-2.3000004.60000010.70000012.80000010.700000270.64800028.8857250.4800000.00
25%11.1625002.37500010.37500018.35000013.475000106.86425019.0317790.0000000.000000
50%27.7000009.95000017.80000025.95000014.700000428.56700031.7259710.0000001.000000
75%53.12500018.00000025.02500032.20000015.7000001498.260000136.95970410.7750001.000000
max196.19000023.60000030.90000038.20000018.5000005672.650000218.496862783.2000001.000000

Highlighting rows that were recorded only in the 2000s and moving it into its own dataframe

The2000s = df.iloc[142:282]
The2000s

DatePPTTminTavTmaxdeltSMIETFlow_cfsFlow?

1422000-019.720.38.516.716.482.62013.3843910.000
1432000-025.823.512.421.417.972.16826.5232470.000
1442000-03105.426.614.722.816.21549.67444.20164151.201
1452000-048.103.618.927.316.7157.62675.2291680.371
1462000-0517.5317.725.633.415.7448.768147.8410502.511
...
2572009-0830.9320.928.235.614.7872.226172.0315970.271
2582009-0957.7216.022.529.013.01286.700100.07444512.701
2592009-1025.809.516.623.614.1428.28052.93313672.401
2602009-110.044.613.322.017.40.53231.0629040.000
2612009-1252.53-2.44.711.814.2246.8914.2318090.041
120 rows x 10 columns

Using the describe function to obtain basic and statistical information of the 2000s dataframe

The2000s.describe()

PPTTminTavTmaxdeltSMIETFlow_cfsFlow?

count109.000000110.000000110.000000110.000000110.000000110.000000110.000000110.000000110.000000
mean46.92800010.65000017.78583324.91916714.269167939.83693380.25104811.2666830.808333
std44.2291438.0523087.7852917.6646991.9304601049.03865262.21598130.1418900.395233
min0.000000-2.4000003.5000008.8000009.9000000.0000002.4530930.0100001.000000
25%10.9350003.10000010.50000018.47500012.875000132.38350021.3895660.0000001.000000
50%36.05000010.90000018.20000025.40000014.700000442.56700066.4394621.0950001.000000
75%66.80000018.32500025.22500031.72500015.7000001498.260000136.95970410.7750001.000000
max205.32000023.40000030.50000037.60000018.2000005317.788000210.456467789.3000001.000000

Highlighting rows to only present data in the 2010s and moving it into its own dataframe

The2010s = df.iloc[282:]
The2010s

DatePPTTminTavTmaxdeltSMIETFlow_cfsFlow?

2822010-0139.64-1.85.612.914.7221.98410.8295500.981
2832010-0275.24-0.55.311.111.6386.7725.2170037.471
2842010-0348.904.112.019.915.8586.80029.9796513.151
2852010-04115.8310.217.424.514.32015.44276.2231680.371
2862010-0575.9214.821.828.914.1565.056108.7159672.771
...
3682018-118.101.89.116.314.573.71015.0391820.101
3692018-1230.350.16.613.012.9200.3108.1022000.011
3702019-016.00-0.16.813.813.940.8008.7402910.000
3712019-023.711.69.417.215.634.87415.6133470.000
3722019-0319.364.310.917.513.2211.02424.9428570.151
111 rows x 10 columns

Using the describe function to obtain basic and statistical information of the 2010s dataframe

The2010s.describe()

PPTTminTavTmaxdeltSMIETFlow_cfsFlow?

count110.000000111.000000111.000000111.000000111.000000111.000000111.000000111.000000111.000000
mean42.73807210.51081817.88226225.19279914.880180885.07253262.09777415.9895900.567566
std45.6755648.5520748.2610138.0684991.9424361074.39555966.41139882.5740020.467660
min0.000000-3.0000004.60000010.4000009.9000000.0000004.3378780.0000000.000000
25%11.9600002.1000009.35000017.15000012.800000144.62250015.4545620.0000001.000000
50%31.10000010.90000016.10000025.90000014.400000405.87000067.7895910.0000001.000000
75%62.31500018.85000025.70000032.05000015.80000013431.070000144.2564130.0000001.000000
max251.24000024.30000031.50000038.80000019.4000005672.650000218.125353709.9000001.000000

Highlighting rows to only present data in the 2010s WITH NO OUTLIERS and moving it into its own dataframe

The2010sNO_OUTLIERS = The2010s[The2010s["Flow?"] < 500]
The2010sNO_OUTLIERS.describe()

PPTTminTavTmaxdeltSMIETFlow_cfsFlow?

count109.000000109.000000109.000000109.000000109.000000109.000000109.000000109.000000109.000000
mean38.99522910.40275217.77339425.14433714.741284118.84263766.95252412.3818050.498724
std36.4846838.5970338.1342078.1342031.0606102.0000000.00000066.95252412.3818050.498724
min0.000000-3.8000004.60000010.40000010.40000010.2000000.0000000.0000000.000000
25%11.7500002.1000009.30000017.10000013.600000143.71500015.1313470.0000001.000000
50%30.35000010.50000018.00000025.80000014.500000457.57000065.1729430.0000001.000000
75%58.95000019.10000025.70000032.20000015.8000001369.377000148.2685962.7700001.000000
max169.51000024.30000031.50000038.80000019.4000005673.206000218.12535372.4000001.000000

Establishing a dataframe of the 1990s that is all non-flow states

NoFlowPeriods90s = The1990s[The1990s["Flow?"] < 0.5]
NoFlowPeriods90s.shape

(46, 10)

Establishing a dataframe of the 2000s that is all non-flow states

NoFlowPeriods00s = The2000s[The2000s["Flow?"] < 0.5]
NoFlowPeriods00s.shape

(23, 10)

Establishing a dataframe of the 2010s that is all non-flow states

NoFlowPeriods10s = The2010s[The2010s["Flow?"] < 0.5]
NoFlowPeriods10s.shape

(48, 10)

Creating the linear model using precipitation values to predict the average streamflow and fitting it

Establishing key values of intercept, slope, R^2, and RMSE

import statsmodels.formula.api as smf

Initialize and fit linear regression model using 'statsmodels'
model = smf.ols("Flow_cfs ~ PPT", data=df) # model object constructor syntax
model = model.fit()

intercept = model.params[0]
slope = model.params[1]
Rsquare = model.rsquared
RMSE = math.sqrt(model.mse_total)

model.params

Intercept-18.335692
PPT0.741168
dtype: float64

Predict values and set up the graph

y_pred = model.predict()

Plot regression against actual data
plt.figure(figsize=(12, 6))
plt.plot(df["PPT"], df["Flow_cfs"], 'o') # scatter plot showing actual data
plt.plot(df["PPT"], y_pred, 'r', linewidth=2) # regression line
plt.xlabel("Recorded Precipitation in Inches")
plt.ylabel("Average Recorded Streamflow")
plt.title("titleline")
plt.show();

Predicted Streamflow from Precipitation Values

R squared = 0.378
RMSE = 51.73

Average Recorded Streamflow

Average Recorded Precipitation in Inches

Creating the linear model using average temperature values to predict the average streamflow and fitting it

Establishing key values of intercept, slope, R^2, and RMSE

import statsmodels.formula.api as smf

Initialize and fit linear regression model using 'statsmodels'
model = smf.ols("Flow_cfs ~ Tav", data=df) # model object constructor syntax
model = model.fit()

intercept = model.params[0]
slope = model.params[1]
Rsquare = model.rsquared
RMSE = math.sqrt(model.mse_total)

model.params

Intercept-1.2168093
Tav0.888591
dtype: float64

Predict values and set up the graph

y_pred = model.predict()

Plot regression against actual data
plt.figure(figsize=(12, 6))
plt.plot(df["Tav"], df["Flow_cfs"], 'o') # scatter plot showing actual data
plt.plot(df["Tav"], y_pred, 'r', linewidth=2) # regression line
plt.xlabel("Average Temperature Degrees Fahrenheit")
plt.ylabel("Average Recorded Streamflow")
plt.title("titleline")
plt.show();

Predicted Streamflow from Average Temperature Values

R squared = 0.015
RMSE = 51.73

Average Recorded Streamflow

Average Temperature Degrees Fahrenheit

Setting up the heatmap to display the correlations

import seaborn as sns
plt.figure(figsize=(15,10))
corr = df.corr(method = 'pearson')
sns.heatmap(corr, annot=True)
plt.show()

PPTTminTavTmaxdeltSMIETFlow_cfsFlow?

PPT10.380.320.25-0.620.940.290.410.45
Tmin0.3810.990.97-0.360.570.970.150.23
Tav0.320.9910.99-0.250.530.980.120.2
Tmax0.250.970.991-0.130.470.960.0920.16
delt-0.62-0.36-0.25-0.131-0.57-0.29-0.28-0.38
SMI0.940.570.530.47-0.5710.50.580.42
ET0.290.970.980.96-0.290.510.010.18
Flow_cfs0.410.150.120.092-0.28-0.580.110.18
Flow?0.450.230.20.16-0.38-0.42-0.180.181

Created a histogram to display the distribution of the precipitation of the 1990s

plt.hist(The1990s["PPT"],bins=25)
plt.xlabel("Precipitation Values")
plt.ylabel("Frequency")
plt.title("Frequency of Precipitation during the 1990s")
plt.show()

Frequency of Precipitation during the 1990s

Frequency

Precipitation Values

Created a histogram to display the distribution of the average flowrate of the 1990s

plt.hist(The1990s["Flow_cfs"],bins=25)
plt.xlabel("Average recorded Flowrate")
plt.ylabel("Frequency")
plt.title("Frequency of Flowrate during the 1990s")
plt.show()

Frequency of Flowrate during the 1990s

Frequency

Average recorded Flowrate

Created a histogram to display the distribution of the precipitation of the 2000s

plt.hist(The2000s["PPT"],bins=25)
plt.xlabel("Precipitation Values")
plt.ylabel("Frequency")
plt.title("Frequency of Precipitation during the 2000s")
plt.show()

Frequency of Precipitation during the 2000s

Frequency

Precipitation Values

Created a histogram to display the distribution of the average flowrate of the 2000s

plt.hist(The2000s["Flow_cfs"],bins=25)
plt.xlabel("Average recorded Flowrate")
plt.ylabel("Frequency")
plt.title("Frequency of Flowrate during the 2000s")
plt.show()

Frequency of Flowrate during the 2000s

Frequency

Average recorded Flowrate

Created a histogram to display the distribution of the precipitation of the 2010s

plt.hist(The2010s["PPT"],bins=25)
plt.xlabel("Precipitation Values")
plt.ylabel("Frequency")
plt.title("Frequency of Precipitation during the 2010s")
plt.show()

Frequency of Precipitation during the 2010s

Frequency

Precipitation Values

Created a histogram to display the distribution of the average flowrate of the 2010s

plt.hist(The2010s["Flow_cfs"],bins=25)
plt.xlabel("Average recorded Flowrate")
plt.ylabel("Frequency")
plt.title("Frequency of Flowrate during the 2010s")
plt.show()

Frequency of Flowrate during the 2010s

Frequency

Average recorded Flowrate

Which parameter (between precipitation and temperature values) can be a better predictor for streamflow at the station of study? why?

Between the two parameters, the better predictor for the streamflow is precipitation, as when fitting a linear data model reveals that the precipitation has a vastly superior r-squared value compared to temperature. Although the r-squared value is 0.378 out of 1, which is considered a weak or low effect size; however, it is considerably better than the temperature's r-squared value of 0.015, which gives the impression of no correlation regarding streamflow output.

Periods with the flowrate of 0 (No-flow periods) can be viewed as indicators of drought. What can you understand from comparing the number of no-flow days in the 90s, 2000s, and 2010s? Is the upstream of the Colorado River becoming more or less prone to drought?

By typing in code to analyse the number of no-flow periods throughout the three decades, it is revealed that the 90s had 46 no-flow periods, the 2000s had 23 no-flow periods, and the 2010s had 48 no-flow periods. In analysis, I was quite surprised to see that the amount of no-flow periods in the 2000s differed greatly compared to the other 2 decades. As a result, I created histograms in order to visualize the data set and realized that the low rate of the 2010s had two outliers that greatly influenced the average flow rate, this had made it appear to be on edge with the two other decades at 15.9; however, once removing the outlier, the average flow rate greatly decreased to 5.06 granting me an assumption that the upstream of the Colorado river is becoming more prone to drought since the correlations of the "Flow?" has a high, positive correlation with average flow rate.

What can you infer from comparing the maximum recorded floods in the 90s, 2000s, and 2010s? Based on the records available in the dataset, would you expect to see more or less extreme floods in the future?

By utilizing histograms to view the frequencies of the average flowrate for the past three decades, one can understand that amount of floodings has significantly declined; however, it can be noted that the amount of extreme floods would be expected to increased as before the 2010s, there was no average flowrate above 300, and the 2010s manages to have two extreme floodings at the 500 and 700 mark.