

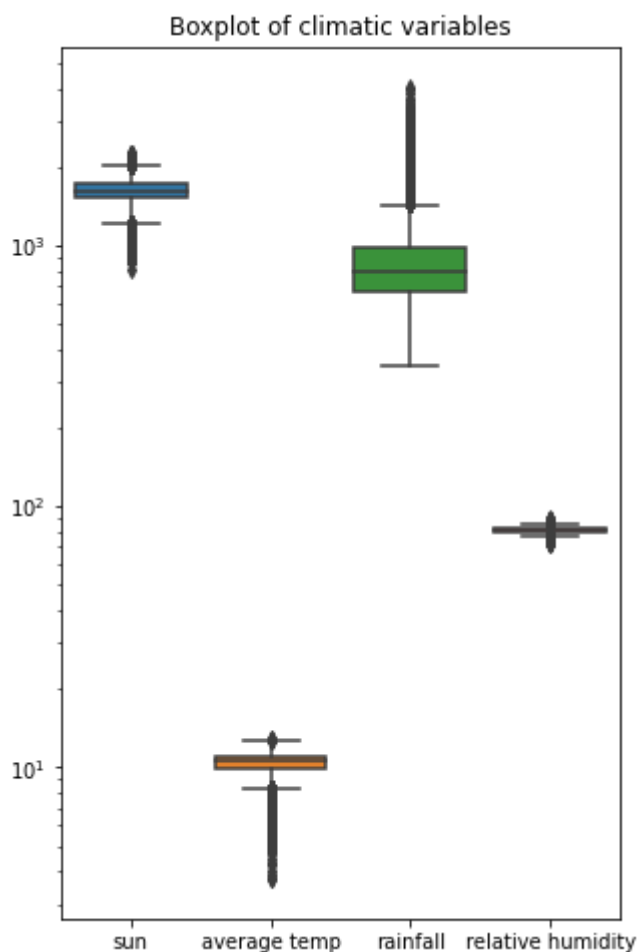
```
In [ ]: import pandas as pd
from sklearn.linear_model import LinearRegression, LogisticRegression
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn import metrics
from sklearn import tree
from sklearn.ensemble import ExtraTreesRegressor
from sklearn.neighbors import KNeighborsClassifier
from sklearn.preprocessing import StandardScaler, MinMaxScaler, RobustScaler
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
from imblearn import under_sampling
```

```
In [ ]: #loading intermediate dataset ukbms and haduk, created in HADUK_intermedia
df = pd.read_csv('BioD_year_site_latlong.csv')
df.drop(columns=df.columns[0], inplace = True)
```

```
In [ ]: df1_climate = df.drop(columns=['Simpsons Index', 'SPECIES RICHNESS', 'TOTAL S
df1_climate.describe()
df2_clim = df.drop(columns=['Site Number', 'YEAR', 'Easting', 'Northing', 'la

fig, ax = plt.subplots(figsize=(5,8))
bxplot = sns.boxplot(data=df1_climate)
bxplot.set_yscale('log')
plt.title('Boxplot of climatic variables')
```

Out[]: Text(0.5, 1.0, 'Boxplot of climatic variables')



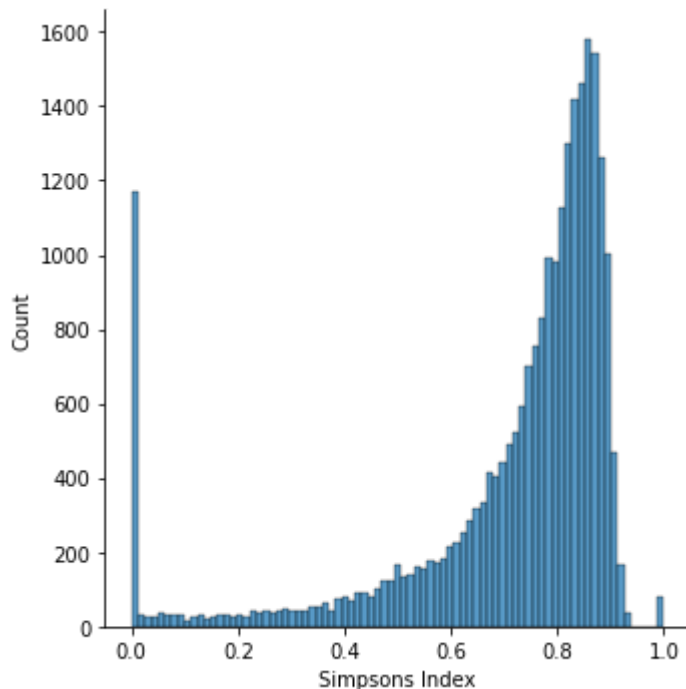
```
In [ ]: df2_clim.describe()
```

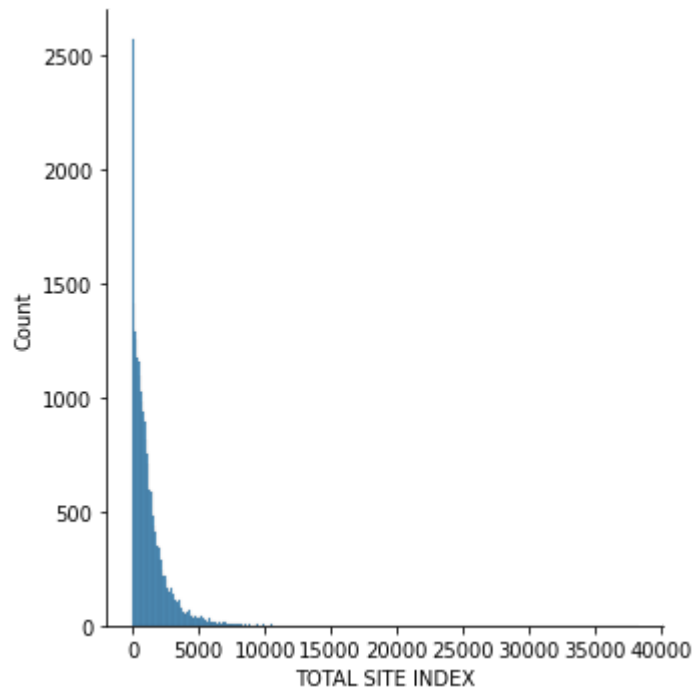
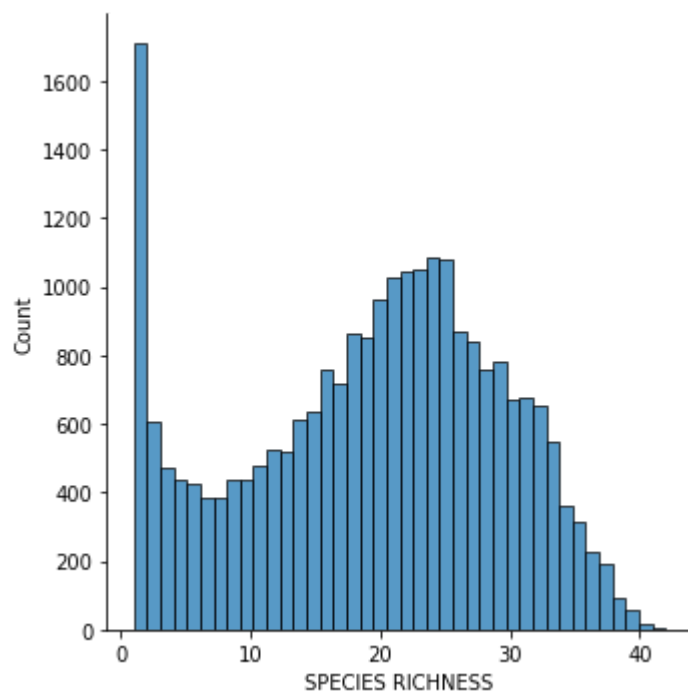
```
Out[ ]:
```

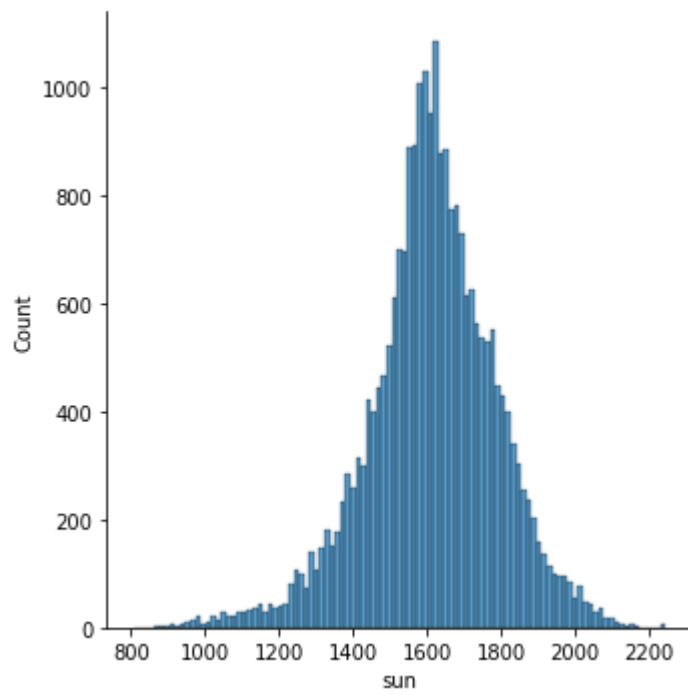
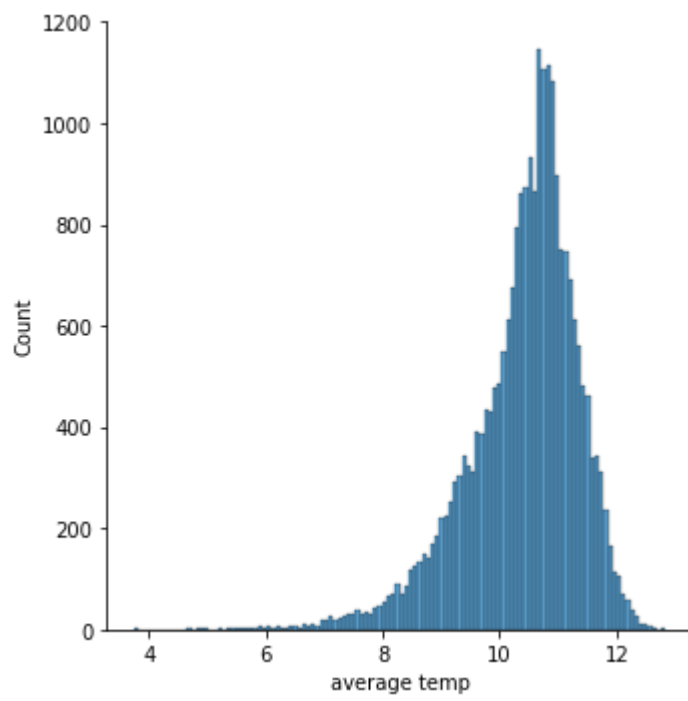
	SPECIES RICHNESS	TOTAL SITE INDEX	sun	average temp	rainfall	rel hurr
count	24558.000000	24558.000000	24558.000000	24558.000000	24558.000000	24558.00
mean	19.234425	1242.418682	1613.894576	10.396600	865.006929	81.34
std	9.719643	1444.173563	174.682971	0.968519	303.500379	1.86
min	1.000000	2.000000	806.658881	3.743084	346.315933	72.32
25%	12.000000	332.000000	1519.217458	9.918850	668.173773	80.3
50%	21.000000	839.000000	1616.506798	10.573544	801.808547	81.41
75%	27.000000	1639.000000	1722.647021	11.021346	976.247667	82.44
max	42.000000	38246.000000	2244.008016	12.839887	4071.312595	90.01

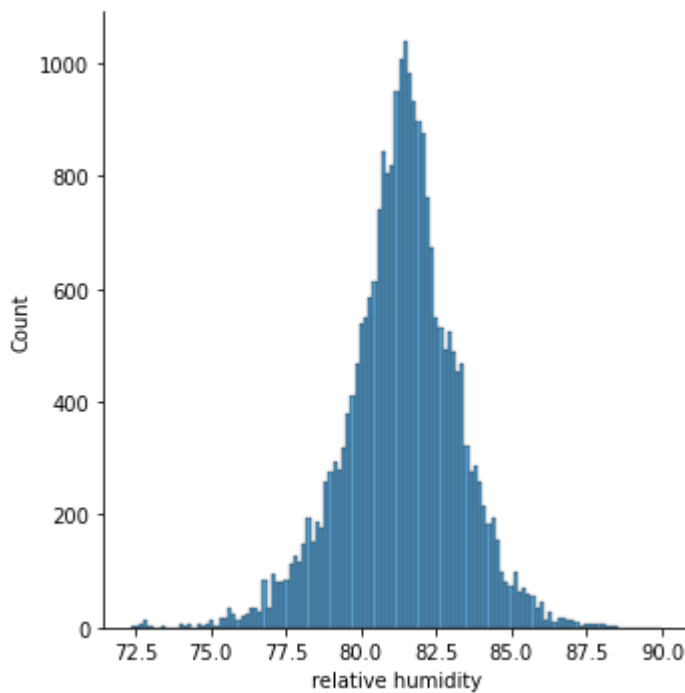
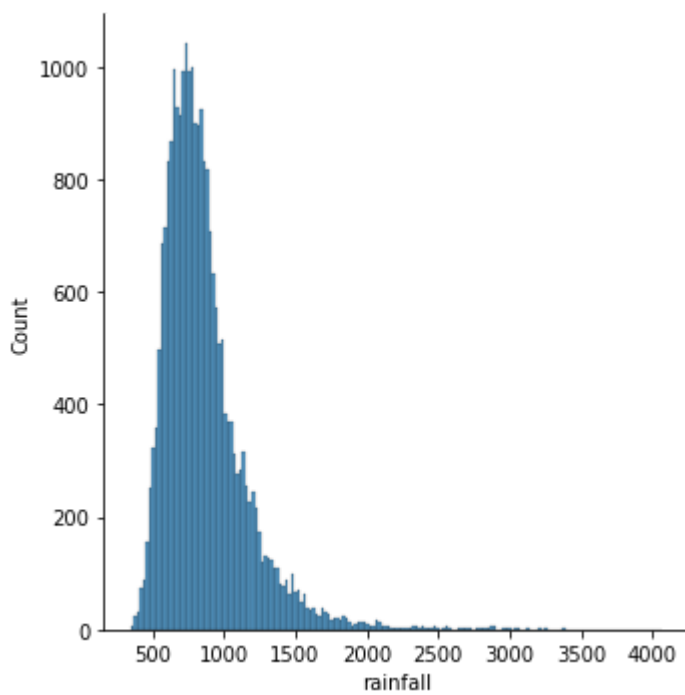
```
In [ ]: sns.displot(df['Simpsons Index'])
sns.displot(df['SPECIES RICHNESS'])
sns.displot(df['TOTAL SITE INDEX'])
sns.displot(df['average temp'])
sns.displot(df['sun'])
sns.displot(df['rainfall'])
sns.displot(df['relative humidity'])
```

```
Out[ ]: <seaborn.axisgrid.FacetGrid at 0x7fc28e0528b0>
```







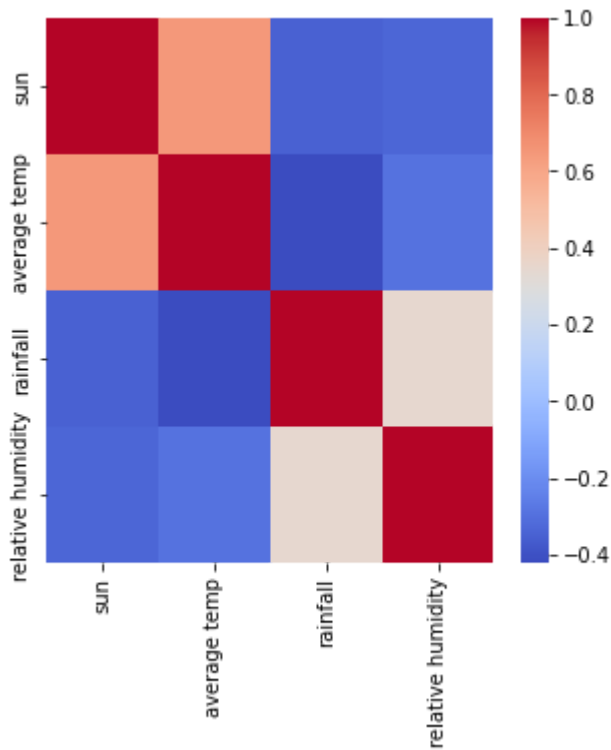


Multiple Linear Regression, carried out on climatic variables (x) and species richness (y):

```
In [ ]: x = df[['sun', 'average temp', 'rainfall', 'relative humidity']]
y = df[['SPECIES RICHNESS']]
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.3, random_
```

```
In [ ]: corr = df1_climate.corr()
plt.figure(figsize=(5,5))
sns.heatmap(corr, cmap='coolwarm')
```

Out[]: <AxesSubplot:>



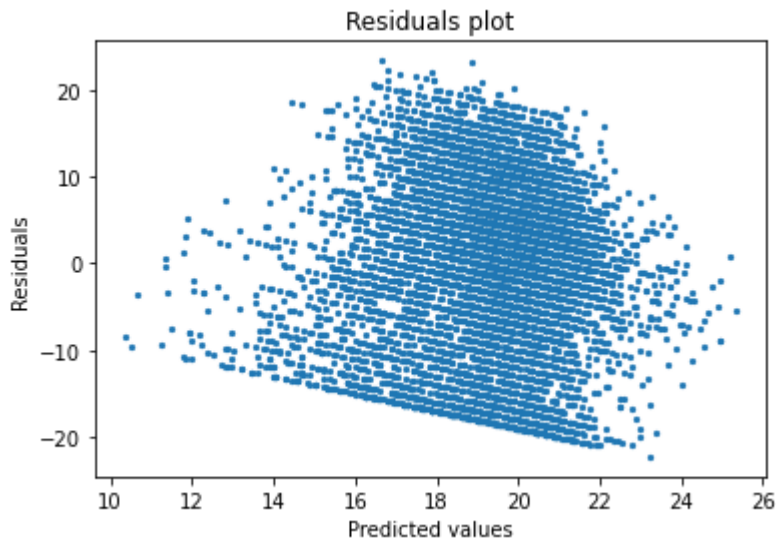
```
In [ ]: x_train.drop(columns=['sun'],inplace=True)
x_test.drop(columns=['sun'],inplace=True)
```

```
In [ ]: #training multiple linear regression model
model = LinearRegression()
model.fit(x_train, y_train)
print(model.coef_)
print(model.intercept_)
```

```
[[ 0.96703309 -0.00211092 -0.45044449]]
[47.625016]
```

```
In [ ]: preds_SR = model.predict(x_test)
```

```
In [ ]: #plotting residuals
residuals = y_test-preds_SR
plt.scatter(preds_SR,residuals,5)
plt.ylabel('Residuals')
plt.xlabel('Predicted values')
plt.title('Residuals plot')
plt.show()
```



```
In [ ]: mse = metrics.mean_squared_error(y_test,preds_SR)
r2 = metrics.r2_score(y_test,preds_SR)
print(mse)
print(np.sqrt(mse))
print(r2)
```

```
91.9585298731896
9.589501023160151
0.039946953528546936
```

```
In [ ]: #Checking diversity
x = df[['average temp','rainfall','relative humidity']]
y= df[['Simpsons Index']]
x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.3, random_
modell = LinearRegression()
modell.fit(x_train, y_train)
print(modell.coef_)
print(modell.intercept_)
preds_Div = modell.predict(x_test)
mse = metrics.mean_squared_error(y_test,preds_Div)
r2 = metrics.r2_score(y_test,preds_Div)
print('mse: ' + str(mse))
print('rmse: ' + str(np.sqrt(mse)))
print('r2: ' + str(r2))
```

```
[[ 9.50157468e-03 -7.94726070e-05 -8.92072927e-03]]
[1.40639564]
mse: 0.04964239506285567
rmse: 0.22280573390928626
r2: 0.03478576404689471
```

Decision tree classifier:

```
In [ ]: #multiclass classifcation decision tree classifier
x_c = df[['sun','average temp','rainfall','relative humidity']]
y_c = df[['Diversity']]
x_ctrain,x_ctest,y_ctrain,y_ctest = train_test_split(x_c,y_c,test_size=0.3,

dtree_model = tree.DecisionTreeClassifier(class_weight='balanced')
dtree_model.fit(x_ctrain,y_ctrain)
dtree_preds = dtree_model.predict(x_ctest)

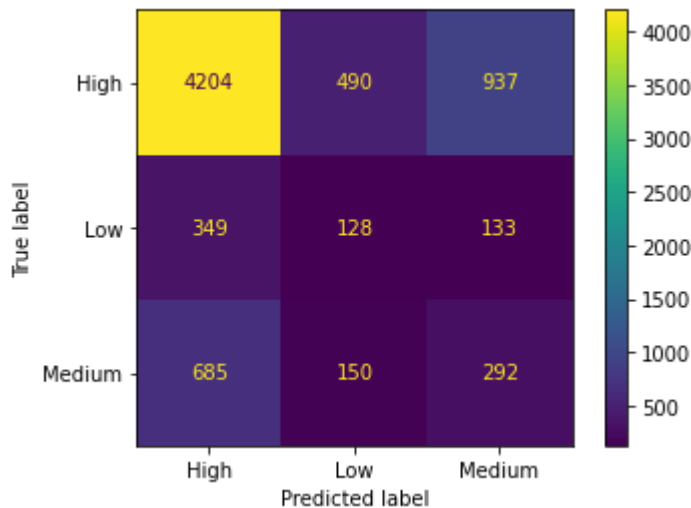
cm = metrics.confusion_matrix(y_ctest,dtree_preds)
acc = metrics.accuracy_score(y_ctest, dtree_preds)
```

```

metrics.ConfusionMatrixDisplay.from_predictions(y_c_test,dtree_preds)
print(metrics.f1_score(y_c_test,dtree_preds, average='weighted'))
print(acc)

```

0.642475852461283
0.6275787187839305



```

In [ ]: #hyperparameter tuning
parameters={"splitter":["best","random"],
            "max_depth" : [1,3,5,7,9,11],
            "min_samples_leaf": [1,3,5,7,9],
            "min_weight_fraction_leaf": [0,0.1,0.2,0.3,0.4],
            "max_features": [1,"log2","sqrt",None],
            "max_leaf_nodes": [None,10,30,50,70,90] }

```

```

In [ ]: #d tree classifier tunin
ctree_tuning = tree.DecisionTreeClassifier(class_weight='balanced')
c_tuning_model=GridSearchCV(ctree_tuning,param_grid=parameters,scoring='acc')
tuning_sample = df.sample(frac=0.3) #sampling dataset to prevent crashes w
x_ctuning = tuning_sample[['sun','average temp','relative humidity','rainfa
y_ctuning = tuning_sample[['Diversity']]

```

```

In [ ]: c_tuning_model.fit(x_ctuning,y_ctuning)

```

```

Out[ ]: ▸ GridSearchCV
        ▸ estimator: DecisionTreeClassifier
          ▸ DecisionTreeClassifier

```

```

In [ ]: print(c_tuning_model.best_params_)
        print(c_tuning_model.best_score_)

{'max_depth': 5, 'max_features': 1, 'max_leaf_nodes': 50, 'min_samples_lea
f': 7, 'min_weight_fraction_leaf': 0, 'splitter': 'random'}
0.7693769611973171

```

```

In [ ]: tuned_ctree = tree.DecisionTreeClassifier(class_weight='balanced',max_depth=
tuned_ctree.fit(x_ctrain,y_ctrain)
tuned_ctree_preds = tuned_ctree.predict(x_c_test)

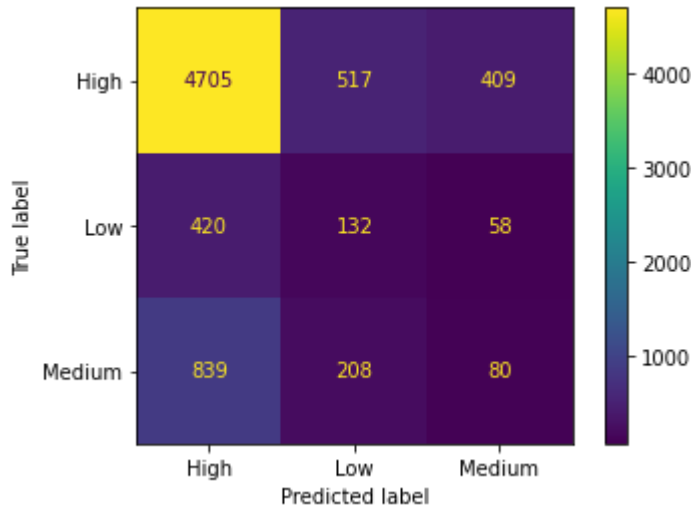
```



```
acc = metrics.accuracy_score(y_test, tuned_ctree_preds)

metrics.ConfusionMatrixDisplay.from_predictions(y_test,tuned_ctree_preds)
print(acc)
```

0.6673452768729642



KNN CLASSIFICATION:

In []:

```
#multiclass classifcation k-nearest neighbours
knn_model = KNeighborsClassifier(n_neighbors = 10).fit(x_train,y_train)
#accuracy = knn_model.score(x_test,y_test)

knn_preds = knn_model.predict(x_test)
cm_knn = metrics.confusion_matrix(y_test,knn_preds)
acc_knn =metrics.accuracy_score(y_test,knn_preds)

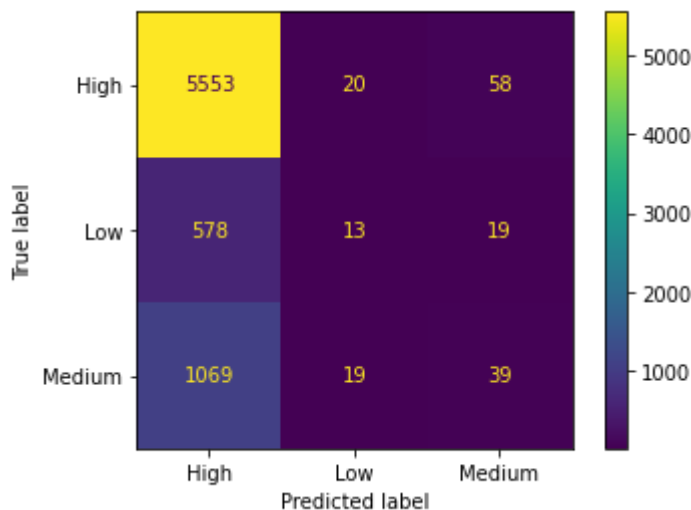
metrics.ConfusionMatrixDisplay.from_predictions(y_test,knn_preds)
print(acc_knn)
metrics.f1_score(y_test,knn_preds, average='weighted')
```

/Library/Frameworks/Python.framework/Versions/3.9/lib/python3.9/site-package
s/sklearn/neighbors/_classification.py:207: DataConversionWarning: A column-
vector y was passed when a 1d array was expected. Please change the shape of
y to (n_samples,), for example using ravel().
return self._fit(X, y)

0.7607220412595005

Out []:

0.6743548782062548



KNN SCALED X VARIABLES:

```
In [ ]: #scaled knn
scaler = StandardScaler()
df_scale = df.copy()
features = [['sun', 'average temp', 'rainfall', 'relative humidity']]
for feature in features:
    df_scale[feature] = scaler.fit_transform(df_scale[feature])

x_scaled = df_scale[['sun', 'average temp', 'rainfall', 'relative humidity']]
y_s = df_scale[['Diversity']]

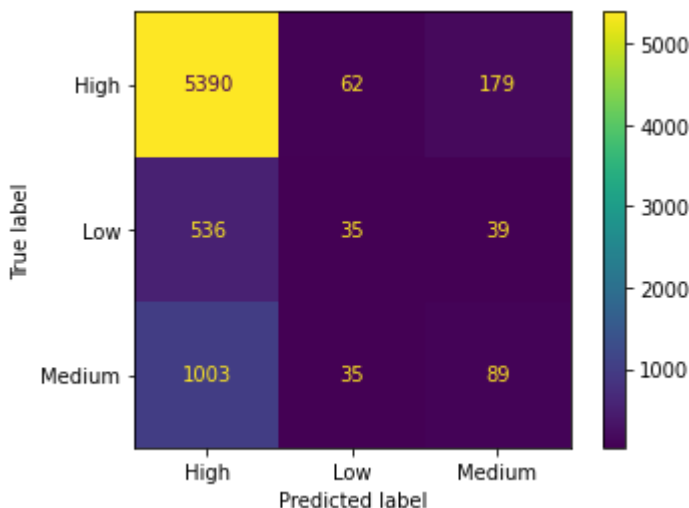
x_strain, x_stest, y_strain, y_stest = train_test_split(x_scaled, y_s, test_size=
knn_scaled_model = KNeighborsClassifier()
knn_scaled_model.fit(x_strain, y_strain)

knn_scaled_preds = knn_scaled_model.predict(x_stest)
acc_scaled = metrics.accuracy_score(y_stest, knn_scaled_preds)

metrics.ConfusionMatrixDisplay.from_predictions(y_stest, knn_scaled_preds)
print(acc_scaled)
```

/Library/Frameworks/Python.framework/Versions/3.9/lib/python3.9/site-packages/sklearn/neighbors/_classification.py:207: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

```
return self._fit(X, y)
0.748371335504886
```



```
In [ ]: metrics.f1_score(y_stest, knn_scaled_preds, average='weighted')
```

```
Out[ ]: 0.682738313321154
```

```
In [ ]: #knn hyperparameter tuning
n_neighbors = [5, 10, 15, 20, 25, 30]
leaf_size = list(range(1, 25))
p = [1, 2]

knn_hyperparameters = dict(n_neighbors=n_neighbors, leaf_size=leaf_size, p=p)
knn_tuned = KNeighborsClassifier()
knn_tuned_model = GridSearchCV(knn_tuned, knn_hyperparameters, cv=3)
```

```
In [ ]: knn_tuned_model.fit(x_strain, y_strain.values.ravel())
```

```
print('best leaf size:', knn_tuned_model.best_estimator_.get_params()['leaf_size'])
print('best p:', knn_tuned_model.best_estimator_.get_params()['p'])
print('best n_neighbors:', knn_tuned_model.best_estimator_.get_params()['n_neighbors'])
```

```
best leaf size: 2
best p: 1
best n_neighbors: 30
```

In []:

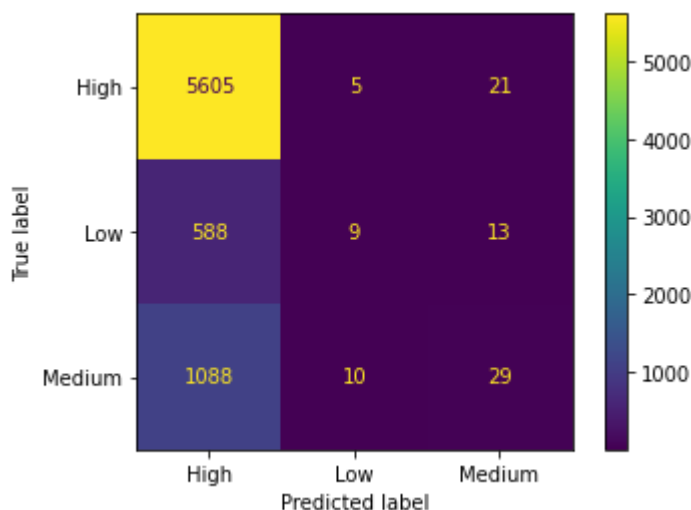
```
tuned_knn_model = KNeighborsClassifier(leaf_size=2,p=1,n_neighbors=30)
tuned_knn_model.fit(x_strain,y_strain)
tuned_knn_preds = tuned_knn_model.predict(x_test)

acc_tuned_knn = metrics.accuracy_score(y_test,tuned_knn_preds)

metrics.ConfusionMatrixDisplay.from_predictions(y_test,tuned_knn_preds)
print(acc_tuned_knn)
print(metrics.f1_score(y_test,tuned_knn_preds, average='weighted'))
```

/Library/Frameworks/Python.framework/Versions/3.9/lib/python3.9/site-packages/sklearn/neighbors/_classification.py:207: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().

```
return self._fit(X, y)
0.7658794788273615
0.6733164631286662
```



Condensed Nearest Neighbor Undersampling:

In []:

```
undersample = under_sampling.CondensedNearestNeighbour()
x_under,y_under = undersample.fit_resample(x_strain,y_strain)
```

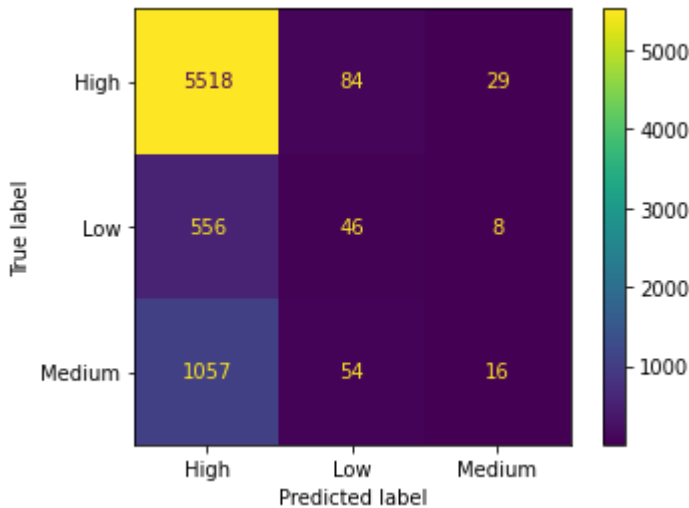
In []:

```
cnn_knn = KNeighborsClassifier(leaf_size=2,p=1,n_neighbors=30)
cnn_knn.fit(x_under,y_under)
cnn_knn_preds = cnn_knn.predict(x_test)

acc_knn_cnn = metrics.accuracy_score(y_test,cnn_knn_preds)
metrics.ConfusionMatrixDisplay.from_predictions(y_test,cnn_knn_preds)
print(acc_knn_cnn)
```

```
/Library/Frameworks/Python.framework/Versions/3.9/lib/python3.9/site-packages/sklearn/neighbors/_classification.py:207: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().
```

```
return self._fit(X, y)
0.757328990228013
```



```
In [ ]: rus_undersample = under_sampling.RandomUnderSampler(random_state=44,replace=False)
x_rus,y_rus = undersample.fit_resample(x_strain,y_strain)
```

```
In [ ]: rus_knn = KNeighborsClassifier(leaf_size=2,p=1,n_neighbors=30)
rus_knn.fit(x_rus,y_rus)
rus_knn_preds = rus_knn.predict(x_test)

acc_knn_rus = metrics.accuracy_score(y_test,rus_knn_preds)
metrics.ConfusionMatrixDisplay.from_predictions(y_test,rus_knn_preds)
print(acc_knn_rus)
```

```
/Library/Frameworks/Python.framework/Versions/3.9/lib/python3.9/site-packages/sklearn/neighbors/_classification.py:207: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().
```

```
return self._fit(X, y)
0.756786102062975
```

