

Lab 5: Data Quality and Validation

Date: Apr 11, 2025

Ryan Varghese, 100870665

# 1. Great Expectations

```
age workclass fillugt education education-num marital-status occupation relationship race sex capital-gain capital-loss hours-per-week native-country income

1 50 Self-emp-not-inc 83311 Bachelors 13 Meried-civ-spouse Exec-managerial Husband White Male 0 0 13 United-States <50K

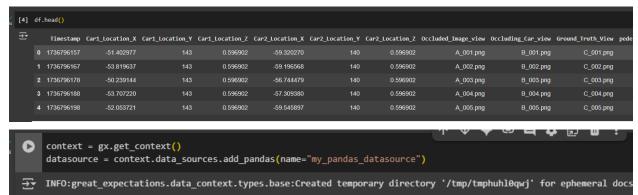
2 58 Private 25466 HS-grad 9 Divorced Handlers-cleaners No-lin-family White Male 0 0 40 United-States <50K

3 53 Private 254721 11th 7 Married-ov-spouse Handlers-cleaners Husband Block Male 0 0 40 United-States <50K

4 28 Private 338409 Bachelors 13 Married-ov-spouse Prof-specially White Black Female 0 0 40 United-States <50K
```

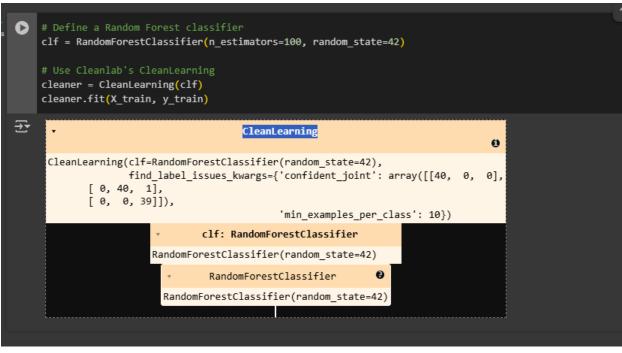
```
validation result = batch.validate(expectation)
print(validation_result)
Calculating Metrics: 100%
                                                                  10/10 [00:00<00:00, 220.28it/s]
 "success": true,
  "expectation_config": {
    "type": "expect_column_values_to_be_between",
      "batch_id": "pandas-pd dataframe asset",
"column": "education-num",
      "min value": 0.0,
      "max_value": 20.0
    "meta": {}
 },
"result": {
    "element_count": 32561,
    "unexpected_count": 0,
    "unexpected_percent": 0.0,
    "partial_unexpected_list": [],
    "missing_count": 0,
    "missing_percent": 0.0,
    "unexpected_percent_total": 0.0,
    "unexpected_percent_nonmissing": 0.0,
    "partial_unexpected_counts": [],
    "partial_unexpected_index_list": []
  "meta": {},
  "exception_info": {
    "raised_exception": false,
    "exception_traceback": null,
    "exception_message": null
```

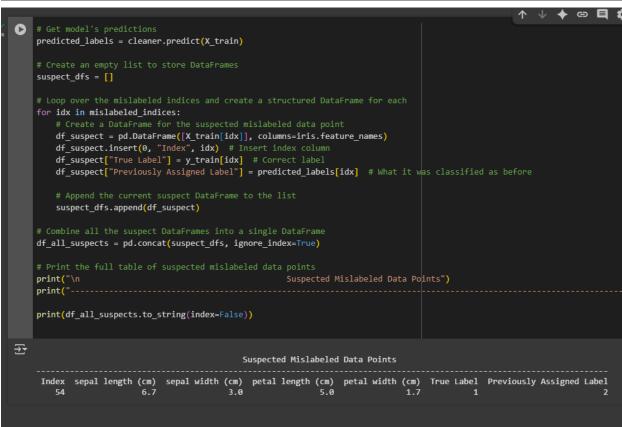
### Task 1:



## 2. CleanLab

#### Task 2:





### Task 3:

[6.

2.2

5.

```
# Load the Iris dataset
    iris = load_iris()
    df = pd.DataFrame(iris.data, columns=iris.feature_names)
    df['target'] = iris.target
    print(df.head())
₹
       sepal length (cm) sepal width (cm) petal length (cm) petal width (cm) \
                    5.1
                                     3.5
                                                       1.4
                                                                        0.2
                    4.9
                                     3.0
                                                       1.4
                                                                        0.2
    1
    2
                    4.7
                                     3.2
                                                       1.3
                                                                        0.2
                    4.6
                                     3.1
                                                       1.5
                                                                        0.2
    4
                    5.0
                                     3.6
                                                       1.4
                                                                        0.2
       target
    0
           0
    1
           0
           0
    2
    3
           0
    4
           0
     X = iris.data
      y = iris.target
      # Use CleanLearning for anomaly detection
      clf = CleanLearning()
      clf.fit(X, y)
      # Find potential anomalies in labels
      label_issues = clf.find_label_issues(X, y)
      # Output the anomalies
      anomalies = np.where(label_issues["is_label_issue"])[0]
      print(f"Anomalies detected at indices: {anomalies}")
      print(f"Suspected anomaly values: {X[anomalies]}")
 → Anomalies detected at indices: [ 18 31 68 82 106 119]
      Suspected anomaly values: [[5.7
                                       3.8
                                                       5.82076585 0.3
      [5.4
                  3.4
                         5.57950291 0.4
                 2.2
                            6.61624076 1.5
       [6.2
       [5.8
                  2.7
                            6.26680751 1.2
       [4.9
                  2.5
                             4.5
                                        1.7
                                        1.5
                                                  ]]
```

```
# Create an empty list to store DataFrames
    suspect_dfs = []
    flower_species = {0.0: "Setosa", 1.0: "Versicolor", 2.0: "Virginica"}
    for idx in anomalies:
        # Create a DataFrame for the suspected anomaly data point
        df_suspect = pd.DataFrame([df.iloc[idx][iris.feature_names].values], columns=iris.feature_names)
        df_suspect.insert(0, "Index", idx) # Insert index column
        df_suspect["True Label"] = df.iloc[idx]["target"]
df_suspect["Flower Species"] = flower_species[y[idx]] # Map label to flower species
        # Append the current suspect DataFrame to the list
        suspect_dfs.append(df_suspect)
    df_all_suspects = pd.concat(suspect_dfs, ignore_index=True)
                                                    Suspected Anomalous Data Points")
    print("---
    print(df_all_suspects.to_string(index=False))
<del>____</del>
                                           Suspected Anomalous Data Points
     Index sepal length (cm) sepal width (cm) petal length (cm) petal width (cm) True Label Flower Species
                                                                                           0.0
                                                                                                        Setosa
                                3.8 5.820766
3.4 5.579503
                                                                      0.3
        18
                          5.4
                                                                                   0.4
                                                                                               0.0
                                                                                                            Setosa
                                                   6.616241
6.266808
4.500000
5.000000
                                                                                                        Versicolor
                                                                                               1.0
                                                                                                       Versicolor
                          5.8
                          4.9
                                                                                               2.0
                                                                                                        Virginica
                                                                                                        Virginica
                           6.0
```

https://github.com/RyanJohnV/RyanVarghese 100870665 SoftQuality Lab5.git