

# 2882. Drop Duplicate Rows

Easy

[Leetcode Link](#)

## Problem Description

The problem presents a DataFrame named `customers` with columns `customer_id`, `name`, and `email`. We are informed that there are duplicate records based on the `email` column. The task is to remove these duplicate rows, but crucially, we must keep only the first occurrence of each `email`. The DataFrame may contain unique `customer_id` values and `name` values, but some `emails` are associated with more than one customer. The goal is to return a DataFrame where all `email` addresses appear only once, preserving the record of the first customer who had that email.

## Intuition

The `drop_duplicates` function from the `pandas` library immediately comes to mind for this task. This function is tailor-made for situations like this, where we need to remove duplicate rows based on one or multiple column values. The function comes with the `subset` parameter, allowing us to specify on which columns to check for duplicates. In this case, we'll set the `subset` to `['email']`, so the function will only consider the `email` column for finding duplicates.

By default, `drop_duplicates` keeps the first occurrence of a duplicated row, which aligns perfectly with our requirements. We do not need to set the `keep` argument as its default value is `'first'`. We can also safely ignore the `inplace` parameter, or set it to `False`, as we want to return a new DataFrame rather than modify the original `customers` DataFrame in place. This provides a simple and efficient solution to the problem that can be implemented and understood with minimal code.

The solution code leverages the described functionality provided by `pandas`, and thus with a single line of code, we achieve the desired output, effectively removing all duplicates based on the `email` while keeping the first occurrence of each.

## Solution Approach

The solution is straightforward due to the robust capabilities of the Pandas library in Python, which is designed to handle and manipulate data in DataFrame structures.

The key steps of the solution are:

1. Import the Pandas library to use its functionalities.
2. Define a function `dropDuplicateEmails`, which takes a `customers` DataFrame as its argument.
3. Use the `drop_duplicates` method available in the Pandas library to remove duplicate rows based on specific column values. The `drop_duplicates` method identifies and removes duplicate rows with the following considerations:
  - The `subset` parameter specifies the columns to consider for identifying duplicate rows. In our case, this is set to `['email']` so that the method looks for duplicates only in the `email` column.
  - By default, the `keep` parameter is set to `'first'`, which means that if duplicates are found, the first occurrence is kept while the subsequent duplicates are removed. This behavior is exactly what we need to solve this problem, so we do not need to specify this parameter explicitly.
  - The method returns a new DataFrame with the duplicates removed, which we immediately return from our function.

Therefore, the core algorithm involves no complex loops or conditionals due to the high-level abstraction provided by Pandas. The data structure used is the DataFrame itself, and the pattern applied is the direct use of a library function designed for this exact purpose.

The implementation is as follows:

```
1 import pandas as pd
2
3 def dropDuplicateEmails(customers: pd.DataFrame) -> pd.DataFrame:
4     return customers.drop_duplicates(subset=['email'])
```

This elegant and concise implementation leverages Pandas' capabilities to solve the problem with minimal code and time complexity.

## Example Walkthrough

Let's walk through a small example to illustrate how the solution approach effectively removes duplicate email addresses and retains only the first occurrence of each.

Suppose we have a small `customers` DataFrame that looks like this:

customer_id	name	email
1	Alice	alice@example.com
2	Bob	bob@example.com
3	Alice	alice@example.com
4	Charlie	charlie@example.com
5	David	bob@example.com

Here, we can see that `Alice` appears twice with the same email (`alice@example.com`), and so does `Bob` with his email (`bob@example.com`). We would like to have only the first occurrence of each unique email address.

Following the solution, this is what we would do:

1. Import the pandas library so that we can work with DataFrames.
2. Create the `dropDuplicateEmails` function.
3. Use the `drop_duplicates` method on the `customers` DataFrame with `subset=['email']`.

When we apply the function `dropDuplicateEmails` to our DataFrame, the `drop_duplicates` method processes the DataFrame as follows:

- It checks the `email` column for duplicate values since we've set `subset=['email']`.
- When it finds the duplicate values `alice@example.com` and `bob@example.com`, it keeps the first occurrence of each (the rows with `customer_id` 1 and 2) and discards the other occurrences (the rows with `customer_id` 3 and 5).

The DataFrame returned by the `dropDuplicateEmails` function will be:

customer_id	name	email
1	Alice	alice@example.com
2	Bob	bob@example.com
4	Charlie	charlie@example.com

So the final DataFrame correctly contains only one record for each email address, and the first customers with those emails have been kept. The duplicate rows have been removed, achieving our desired result. The elegance of this solution lies in its simplicity and the use of Pandas' high-level functionality, which makes the code clean, readable, and efficient.

## Python Solution

```
1 import pandas as pd # Importing the pandas library
2
3 def dropDuplicateEmails(customers: pd.DataFrame) -> pd.DataFrame:
4     """
5     Remove duplicate rows from the customers DataFrame based on the 'email' column.
6
7     Parameters:
8     customers (pd.DataFrame): DataFrame containing customer data with an 'email' column.
9
10    Returns:
11    pd.DataFrame: A new DataFrame without duplicate emails.
12    """
13    # Use the drop_duplicates method on the 'customers' dataframe,
14    # specifying 'email' as the subset to identify duplicates by the 'email' column only.
15    unique_customers = customers.drop_duplicates(subset='email')
16
17    # Return the resulting DataFrame with duplicates removed.
18    return unique_customers
19
```

## Java Solution

```
1 import java.util.HashMap;
2 import java.util.Map;
3 import java.util.List;
4 import java.util.ArrayList;
5 import java.util.Objects;
6
7 class Customer {
8     String email;
9     // Other customer fields can be added here
10
11     // Constructor
12     public Customer(String email) {
13         this.email = email;
14     }
15
16     // Email getter
17     public String getEmail() {
18         return email;
19     }
20
21     // Equals and hashCode methods are overridden to use 'email' field for comparison
22     // and to ensure consistent behavior when used as a key in a HashMap.
23     @Override
24     public boolean equals(Object o) {
25         if (this == o) return true;
26         if (!(o instanceof Customer)) return false;
27         Customer customer = (Customer) o;
28         return Objects.equals(email, customer.email);
29     }
30
31     @Override
32     public int hashCode() {
33         return Objects.hash(email);
34     }
35
36     // A toString method for easy printing of customer information.
37     @Override
38     public String toString() {
39         return "Customer{" +
40             "email='" + email + '\'' +
41             '}';
42     }
43 }
44
45 public class DuplicateEmailsRemover {
46
47     /**
48     * Remove duplicate rows from the list of customers based on the 'email' field.
49     *
50     * @param customers (List<Customer>): A list containing customer objects.
51     * @return List<Customer>: A new list without duplicate emails.
52     */
53     public static List<Customer> dropDuplicateEmails(List<Customer> customers) {
54         // Use a HashMap to track unique emails.
55         Map<String, Customer> uniqueCustomersMap = new HashMap<>();
56
57         // Iterate over the list of customer objects
58         for (Customer customer : customers) {
59             // If the email has not been seen before, add the customer to the map.
60             if(!uniqueCustomersMap.containsKey(customer.getEmail())) {
61                 uniqueCustomersMap.put(customer.getEmail(), customer);
62             }
63         }
64
65         // Return the unique customers as a new list (values of the map).
66         return new ArrayList<>(uniqueCustomersMap.values());
67     }
68
69     // A main method to test the functionality
70     public static void main(String[] args) {
71         // Create a list of customers with some duplicate emails
72         List<Customer> customers = new ArrayList<>();
73         customers.add(new Customer("alice@example.com"));
74         customers.add(new Customer("bob@example.com"));
75         customers.add(new Customer("charlie@example.com"));
76         customers.add(new Customer("alice@example.com")); // Duplicate
77
78         // Remove duplicates
79         List<Customer> uniqueCustomers = dropDuplicateEmails(customers);
80
81         // Print out the unique customer list
82         uniqueCustomers.forEach(System.out::println);
83     }
84 }
85
```

## C++ Solution

```
1 #include <iostream>
2 #include <vector>
3 #include <string>
4 #include <algorithm> // For std::unique and std::stable_sort
5 #include "DataFrame.h" // Assuming DataFrame is a fictional class that needs to be included
6
7 // Define a custom comparison operator
8 bool compareByEmail(const Customer& first, const Customer& second) {
9     return first.getEmail() < second.getEmail();
10 }
11
12 // Define a custom equality operator
13 bool equalByEmail(const Customer& first, const Customer& second) {
14     return first.getEmail() == second.getEmail();
15 }
16
17 // Remove duplicate rows from the customers DataFrame based on the 'email' column.
18 DataFrame dropDuplicateEmails(DataFrame& customers) {
19     // Assuming 'customers' is a DataFrame containing 'Customer' objects
20     // and 'Customer' has a method 'getEmail()' to access the 'email' attribute.
21
22     // First, sort the customers by email using the custom comparison operator
23     std::stable_sort(customers.begin(), customers.end(), compareByEmail);
24
25     // Then, use the unique algorithm with a custom equality operator
26     // unique algorithm will require the duplicates to be next to each other
27     // hence the need for the sort operation beforehand.
28     auto lastUnique = std::unique(customers.begin(), customers.end(), equalByEmail);
29
30     // Erase the non-unique elements from the container
31     customers.erase(lastUnique, customers.end());
32
33     // Return the resulting DataFrame with duplicates removed.
34     return customers;
35 }
36
```

## Typescript Solution

```
1 interface Customer {
2     [key: string]: any; // An index signature to allow any string as a key and its value can be anything
3     email: string;      // Ensure that 'email' is always a string
4 }
5
6 /**
7  * Remove duplicate rows from the customers array based on the 'email' property.
8  *
9  * @param customers - Array containing customer objects with an 'email' property.
10  * @returns A new array without duplicate emails.
11  */
12 function dropDuplicateEmails(customers: Customer[]): Customer[] {
13     const seenEmails = new Set<string>(); // To track already encountered emails
14
15     // Use filter to exclude duplicates
16     const uniqueCustomers = customers.filter(customer => {
17         // If seenEmails does not have this email, add it and keep the customer
18         if (!seenEmails.has(customer.email)) {
19             seenEmails.add(customer.email);
20             return true;
21         }
22     });
23
24     // Otherwise, it's a duplicate; exclude it by returning false
25     return uniqueCustomers;
26 }
27
28 // Return the filtered array without duplicate emails
29 return uniqueCustomers;
30
```

## Time and Space Complexity

The time complexity of the `dropDuplicateEmails` function primarily depends on the `drop_duplicates` method of the pandas DataFrame, which in turn depends on the size of the input DataFrame.

### Time Complexity

For pandas `drop_duplicates` method, the time complexity is generally  $O(n)$ , where  $n$  is the number of rows in the DataFrame because it needs to process each row to check for duplicates.

So, the time complexity for the `dropDuplicateEmails` function is also  $O(n)$ .

### Space Complexity

The space complexity for `drop_duplicates` includes the space required to hold the DataFrame and the temporary data structures used to identify duplicates. Typically, it is also  $O(n)$  since a new DataFrame is constructed to store the result without duplicates.

Thus, the space complexity of `dropDuplicateEmails` would also be  $O(n)$ .