Aim:Identifying and handling duplicates using distinct() (R studio ).
Output:





Hariprasad Vishwakarma

S126

SYCS