# 1) Data Preprocessing

*19MID0020*

## Data-Set Creation using Editor



```
      C:\github\winter_semester-2022\CSI3010 Data Warehousing Data Mining\lab exercises\11_1_21\weather.arff - Sublime Text (UNREGISTERED)
   File  Edit  Selection  Find  View  Goto  Tools  Project Preferences  Help

        weather.arff              ✕

   1    @relation weather
   2    @attribute outlook {sunny,rainy,overcast}
   3    @attribute temparature numeric
   4    @attribute humidity numeric
   5    @attribute windy {true,false}
   6    @attribute play {yes,no}
   7
   8    @data
   9    sunny,85.0,85.0,false,no
  10    overcast,80.0,90.0,true,no
  11    sunny,83.0,86.0,false,yes
  12    rainy,70.0,86.0,false,yes
  13    rainy,68.0,80.0,false,yes
  14    rainy,65.0,70.0,true,no
  15    overcast,64.0,65.0,false,yes
  16    sunny,72.0,95.0,true,no
  17    sunny,69.0,70.0,false,yes
  18    rainy,75.0,80.0,false,yes
```
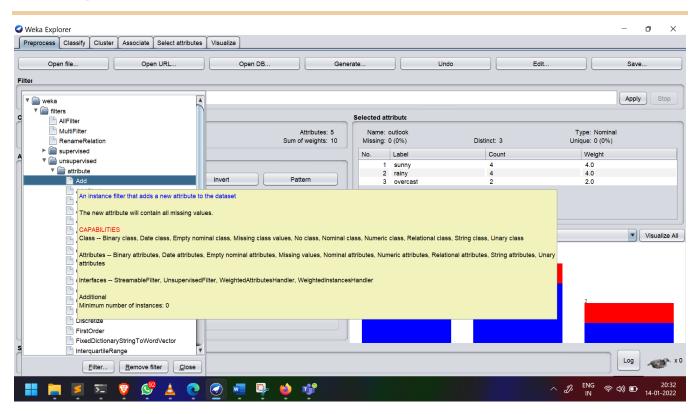
## Opening the created Data-Set using Weka

# Adding new Instances



# Adding Attributes

# Removing Instances

# Removing Attributes

# Normalize

# Replacing the Missing Values

Since there are no missing values in the data-set, I am deleting some records intentionally.

# Remove the values filter