**Figure 1. Apache Nifi Pipeline**

This pipeline was the primary pipeline used to extract, clean, validate, and then load all 4 CSV files. It is only possible to have a pipeline as simple as this one due to the CSV structure matching the structure of our database warehouse. We also chose to not use an intermediary database to load from for this very reason – the intermediary database would have had the exact same structure as our DBW. As such, we chose not to.



**Figure 2. Jupyter Notebook Pandas Cleaning Pipeline**

As mentioned in our report, our machines were unable to handle the size of appointments, as we were unable to properly extract it in Nifi without crashing. As such, we were forced to use alternative means to clean the files with unusable data and split them into more easily processable files. Namely, both px and appointments were filled with issues, such as duplicate IDs and appointments containing vast amounts of entries that did not contain a corresponding px_id. This would cause a foreign key exception error in Nifi, as such, we simply deleted all entries that did not have a corresponding px_id. After all cleaning, data was extracted, cleaned, and loaded in the Nifi pipeline pictured in Figure 1.