



<http://tinyurl.com/2dzbea59>

2 - 3:15PM



<http://tinyurl.com/46czd2ex>

5 - 6:15PM

Bellevue Almshouse dataset

- Data cleaning I
 - Converting data types
 - Duplicates and missing data
 - Frequency of values
- Renaming, deleting, and sorting columns
- Filtering/subsetting data
- Data visualization
- Exporting the dataset

Tidy data structure

- Each **variable** is in a **column**
- Each **observation** is a **row**
- Each **value** is a **cell**

first_name	last_name	age
Mary	Gallagher	28.0
John	Sanin(?)	19.0
Anthony	Clark	60.0
Lawrence	Feeney	32.0

Converting data types

- Converting to date-time data type
 - `data_frame['column_name'] = pd.to_datetime(data_frame['column_name'], format='%Y%m%d')`
 - e.g. `bellevue_df['date_in'] = pd.to_datetime(bellevue_df['date_in'], format='%Y-%m%d')`
- Another way to check data types
 - `data_frame.dtypes`

Dealing with duplicates

- `.duplicated(keep = 'first'/'last'/False):`
 - Creates a True/False dataframe to check which rows in the original dataframe are duplicated
 - `keep`
 - `first`: considers the first entry in the dataframe as the unique entry
 - `last`: considers the last entry in the dataframe as the unique entry
 - `False`: considers all entry as duplicates
 - Default argument: `keep = 'first'`

Dealing with duplicates

- `df[df.duplicated(keep=False)]`
 - Selects duplicated rows from the original dataframe that fulfills the True/False dataframe conditions
- `.drop_duplicates(keep = 'first'/'last'/False):`
 - Drops all the duplicated rows and keeps the first entry, last entry, or none of the entries
 - Default argument: `keep = 'first'`

Missing Data

- `.isna() / .notna()`
 - Creates True/False table for values with/out NA
 - `dataframe_variable['column name'].notna()`
 - E.g. `bellevue_df['professions'].notna()`
 - Filters out NA values by comparing to original df
 - `dataframe_variable[dataframe_variable['column name'].notna()]`
 - E.g. `bellevue_df[bellevue_df['professions'].notna()]`

Missing Data

- `.count()`
 - `count()` method always excludes NaN values
 - To find the percentage of not blank data in every column:
 - `bellevue_df.count() / len(bellevue_df)`
- `.fillna()`
 - Fill the NaN values in the DataFrame with a different value by using the `.fillna()` method
 - `bellevue_df['professions'].fillna('no profession information recorded')`

Frequency: Most common items in a column

- `df["column_name"].value_counts()`
 - To count the number of unique values in a column

Rename Columns

- `.rename(columns={})`
 - `bellevue_df.rename(columns={'professions': 'jobs'})`
 - To save the new column name to the dataframe, we need to overwrite the variable
 - `bellevue_df = bellevue_df.rename(columns={'professions': 'jobs'})`

Adding Columns

- New columns are added to the end of the original data frame
- `Data_frame['new_column'] = <information you plan to store as values in the new column>`
 - Adding columns together
 - `bellevue_df['full_name'] = bellevue_df['first_name'] + bellevue_df['last_name']`
 - Using a list comprehension, I can use conditional statements
 - `bellevue_df['woman'] = ['yes' if gender == 'w' else 'no' for gender in bellevue_df['gender']]`

Drop Columns

- `.drop(columns="column name")`
 - `bellevue_df = bellevue_df.drop(columns="children")`

Sorting Columns

- `.sort_values(by='column_name')`
 - `bellevue_df.sort_values(by='date_in', ascending=True)`

Filter/Subset Data

- `data_frame['column_name'] == 'value'`
 - Produces a True/False table based on condition
 - e.g. `bellevue_df['profession'] == 'teacher'`
- `data_frame[data_frame['column_name'] == 'value']`
 - Filters out the rows from the original data frame that fits the condition
 - e.g. `bellevue_df[bellevue_df['profession'] == 'teacher']`

Groupby Columns

Allows us to group data and perform calculations on the groups

- Creates a groupby object
 - `data_frame.groupby('column_name')`
 - `bellevue_df.groupby('professions')`

Groupby Columns

- Counting non-blank values in each column
 - `data_frame.groupby('column_name').count()`
 - `bellevue_df.groupby('profession').count()`
- Isolating specific column
 - `data_frame.groupby('column_name')['column2'].count()`
 - `bellevue_df.groupby('profession')['gender'].count()`
- Stacking methods
 - `data_frame.groupby('column_name').count().sort_values()`
 - `bellevue_df.groupby('profession')['disease'].count().sort_values(ascending=False)`

Data Visualization

Types of visualization available in pandas

- 'bar' or 'barh' for bar plots
- 'hist' for histogram
- 'box' for boxplot
- 'kde' or 'density' for density plots
- 'area' for area plots
- 'scatter' for scatter plots
- 'hexbin' for hexagonal bin plots
- 'pie' for pie plots

Data Visualization

- `.plot(kind= “<type of visualization>”, title= “<title of graph>”)`
 - Plotting the top 5 disease from Bellevue Almshouse dataset as a bar graph and a title for the graph
 - `bellevue_df['disease'].value_counts()[:5].plot(kind='bar', title='Bellevue Almshouse: Most Frequent "Diseases"')`

Exporting the dataset

- `dataframe.to_csv("<name of file.csv>", encoding= 'utf-8')`
 - `bellevue_df.to_csv("bellevue_df_clean.csv", encoding='utf-8', index=False)`
 - `index=False` will export the dataframe without python's indexing
 - `encoding= 'utf-8'` will tell python how to encode the file; refer to [Aditya Mukerjee: I Can Text You A Pile of Poo, But I Can't Write My Name](#) for conversation on encoding and its importance to programming