

1 Start coding or [generate](#) with AI.

## ✓ Homework 1 - Ally Hayden

### ✓ Task 1

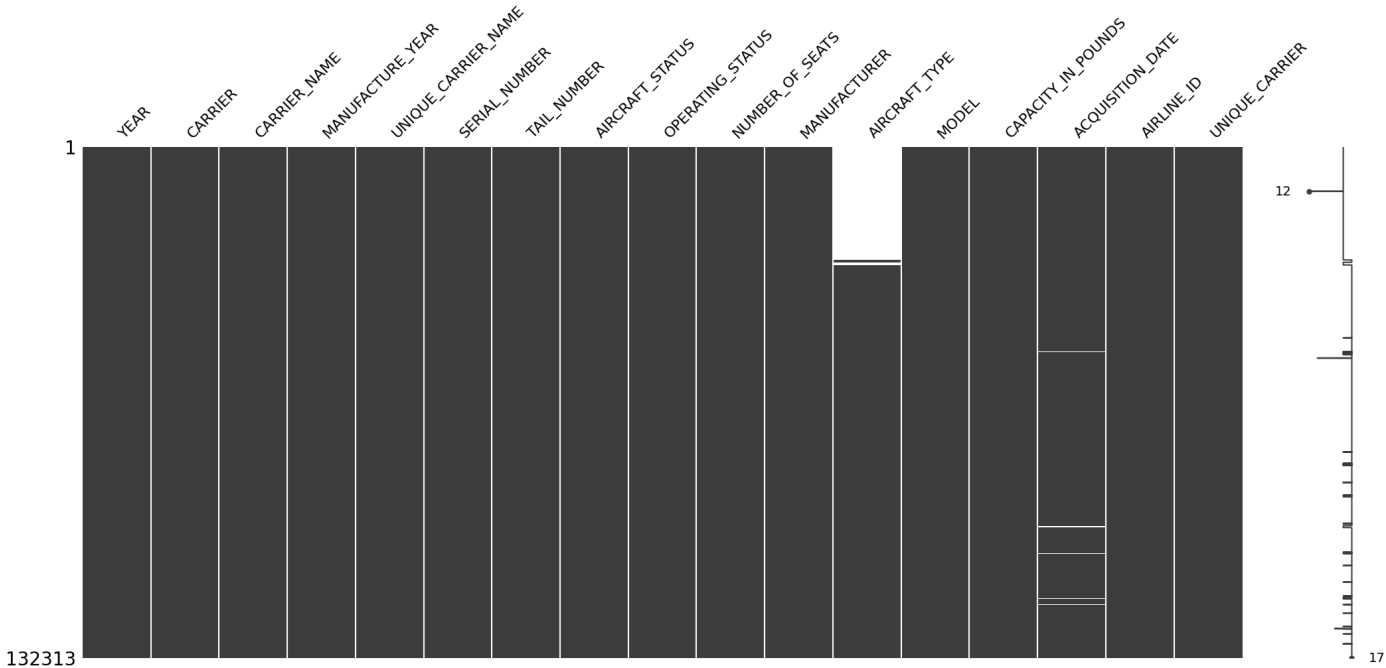
Investigate missing data in aircraft inventory dataset.

```
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4 from google.colab import drive
5 import numpy as np
6
7 # Mount Google Drive
8 drive.mount('/content/drive')
9 pd.set_option('display.precision', 3)
10
11 df = pd.read_csv("/content/drive/MyDrive/DATA 300/T_F41SCHEDULE_B43.csv", keep_default_na=False, na_values=[""])

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True)
<ipython-input-59-beb8561f7e4c>:11: DtypeWarning: Columns (11) have mixed types. Specify dtype option on import or set low_m
df = pd.read_csv("/content/drive/MyDrive/DATA 300/T_F41SCHEDULE_B43.csv", keep_default_na=False, na_values=[""])

1 import missingno as msno
2
3 msno.matrix(df)
```

<Axes: >



```
1 # columns to investigate
2 cols_to_check = ['CARRIER', 'CARRIER_NAME', 'MANUFACTURE_YEAR',
3                 'NUMBER_OF_SEATS', 'CAPACITY_IN_POUNDS', 'AIRLINE_ID']
4
5 # summary
6 missing_summary = df[cols_to_check].isnull().sum().to_frame(name='Missing Count')
7 missing_summary['Missing %'] = (missing_summary['Missing Count'] / len(df)) * 100
8 print("Missing Data Summary:\n", missing_summary)
```

Missing Data Summary:

	Missing Count	Missing %
CARRIER	0	0.000
CARRIER_NAME	105	0.079
MANUFACTURE_YEAR	3	0.002
NUMBER_OF_SEATS	7	0.005
CAPACITY_IN_POUNDS	101	0.076
AIRLINE_ID	105	0.079

CARRIER column

After changing the `pd.read_csv` line to only count blank cells as missing the carrier column has 0 missing values so we do not need to impute it.

CARRIER\_NAME column

```
1 # see if carrier_name and unique_care_name are missing in the same place
2 missing_carrier_name = df[df['CARRIER_NAME'].isna()]
3 print(missing_carrier_name[['CARRIER', 'UNIQUE_CARRIER_NAME']].drop_duplicates())
4
```

```
↗
   CARRIER UNIQUE_CARRIER_NAME
11465      L4                 NaN
54610      OH                 NaN
```

These rows have missing values in both CARRIER\_NAME and UNIQUE\_CARRIER\_NAME. Therefore, they cannot be imputed from UNIQUE\_CARRIER\_NAME. I believe it is safer to leave these rows unimputed because there are so few.

AIRLINE\_ID Column

```
1 # check how many unique AIRLINE_IDs exist
2 carrier_id_check = df.groupby('CARRIER')['AIRLINE_ID'].nunique().sort_values(ascending=False)
3 print(carrier_id_check.head(10)) # shows if some carriers always map to one ID
4
```

```
↗
CARRIER
PT      2
OH      2
YX      2
0WQ     1
16      1
1BQ     1
1EQ     1
16      1
20Q     1
23Q     1
Name: AIRLINE_ID, dtype: int64
```

Some CARRIERS map to only one AIRLINE\_ID:

0WQ, 16, 1BQ, 1EQ, 20Q, 23Q

It is safe to impute these but ONLY these.

```
1 # put it the carriers that map to only ONE airline_id
2 unique_mappings = df.groupby('CARRIER')['AIRLINE_ID'].nunique()
3 single_id_carriers = unique_mappings[unique_mappings == 1].index
4
5 carrier_to_id = df[df['CARRIER'].isin(single_id_carriers)].dropna(subset=['AIRLINE_ID']).drop_duplicates('CARRIER')[['CARRIE
6
7 df['AIRLINE_ID'] = df.apply(
8     lambda row: carrier_to_id[row['CARRIER']] if pd.isna(row['AIRLINE_ID']) and row['CARRIER'] in carrier_to_id else row['AI
9     axis=1
10 )
11
```

CAPACITY\_IN\_POUNDS

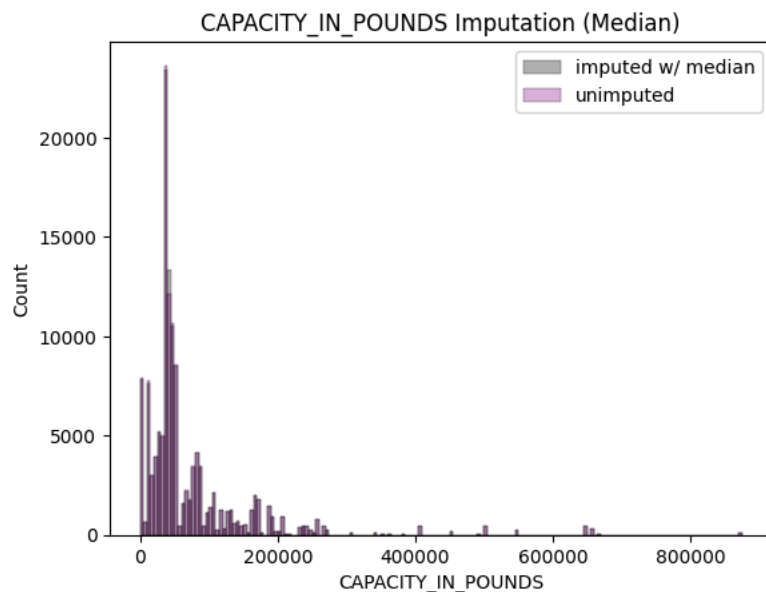
We can impute CAPACITY\_IN\_POUNDS because it's a numeric continuous variable and the data is structured.

I chose Median because it is less influenced by extreme values, KNN on its own operated the same as Median but KNN with features captures inter-variable patterns.

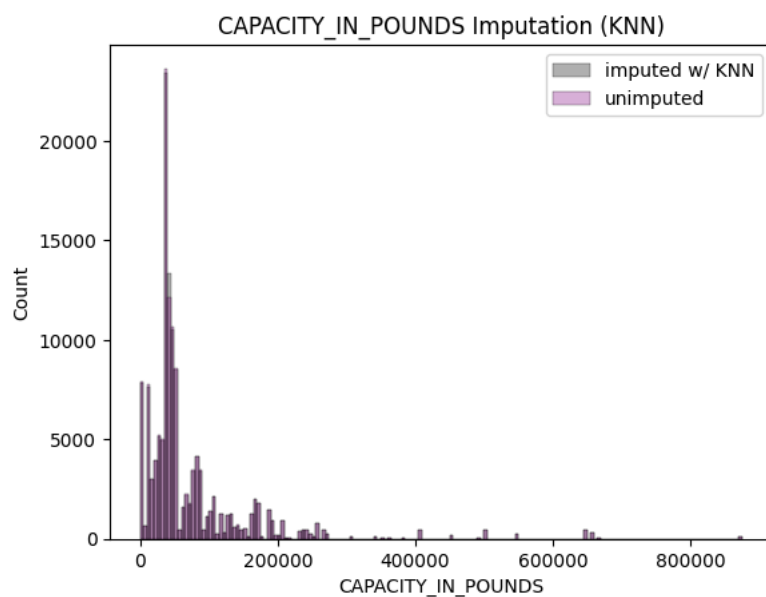
```
1 from sklearn.impute import SimpleImputer, KNNImputer
2 from sklearn.metrics import mean_squared_error
3
4 # random missing
5 df_temp = df[['CAPACITY_IN_POUNDS']].dropna().copy()
6 df_temp.loc[df_temp.sample(frac=0.01, random_state=0).index, 'CAPACITY_IN_POUNDS'] = np.nan
7
8 true_vals_cap = df[['CAPACITY_IN_POUNDS']].dropna().loc[df_temp[df_temp['CAPACITY_IN_POUNDS'].isna()].index]
9
10 # Median
```

```
11 median_imputer_cap = SimpleImputer(strategy='median')
12 df_temp['CAPACITY_IN_POUNDS'] = median_imputer_cap.fit_transform(df_temp[['CAPACITY_IN_POUNDS']])
13
14 # RMSE
15 rmse_cap = np.sqrt(mean_squared_error(true_vals_cap, df_temp.loc[true_vals_cap.index]))
16 print(f"RMSE (Median - CAPACITY_IN_POUNDS): {rmse_cap:.2f}")
17
18 # Plot
19 fig, ax = plt.subplots(1, 1)
20 sns.histplot(df_temp['CAPACITY_IN_POUNDS'], binwidth=5000, alpha=0.3, color='k', label='imputed w/ median')
21 sns.histplot(df['CAPACITY_IN_POUNDS'], binwidth=5000, alpha=0.3, color='purple', label='unimputed')
22 ax.set_title('CAPACITY_IN_POUNDS Imputation (Median)')
23 ax.legend()
24 plt.show()
25
26
27 # KNN
28 knn_imputer_cap = KNNImputer(n_neighbors=5)
29 df_temp['CAPACITY_IN_POUNDS'] = knn_imputer_cap.fit_transform(df_temp[['CAPACITY_IN_POUNDS']])
30
31 # RMSE
32 rmse_cap_knn = np.sqrt(mean_squared_error(true_vals_cap, df_temp.loc[true_vals_cap.index]))
33 print(f"RMSE (KNN - CAPACITY_IN_POUNDS): {rmse_cap_knn:.2f}")
34
35 # Plot
36 fig, ax = plt.subplots(1, 1)
37 sns.histplot(df_temp['CAPACITY_IN_POUNDS'], binwidth=5000, alpha=0.3, color='k', label='imputed w/ KNN')
38 sns.histplot(df['CAPACITY_IN_POUNDS'], binwidth=5000, alpha=0.3, color='purple', label='unimputed')
39 ax.set_title('CAPACITY_IN_POUNDS Imputation (KNN)')
40 ax.legend()
41 plt.show()
42
43 # KNN with features
44 knn_df = df[['CAPACITY_IN_POUNDS', 'NUMBER_OF_SEATS', 'MANUFACTURE_YEAR']].dropna().reset_index(drop=True)
45
46 missing_idx = knn_df.sample(frac=0.01, random_state=0).index
47 true_vals = knn_df.loc[missing_idx, 'CAPACITY_IN_POUNDS'].copy()
48 knn_df.loc[missing_idx, 'CAPACITY_IN_POUNDS'] = np.nan
49
50 knn_imputer = KNNImputer(n_neighbors=5)
51 knn_imputed = pd.DataFrame(knn_imputer.fit_transform(knn_df), columns=knn_df.columns)
52
53 # RMSE
54 rmse_better = np.sqrt(mean_squared_error(true_vals, knn_imputed.loc[missing_idx, 'CAPACITY_IN_POUNDS']))
55 print(f"RMSE (KNN with features - CAPACITY_IN_POUNDS): {rmse_better:.2f}")
56
57 # Plot
58 fig, ax = plt.subplots(1, 1)
59 sns.histplot(knn_df['CAPACITY_IN_POUNDS'], binwidth=5000, alpha=0.3, color='k', label='imputed w/ KNN with features')
60 sns.histplot(df['CAPACITY_IN_POUNDS'], binwidth=5000, alpha=0.3, color='purple', label='unimputed')
61 ax.set_title('CAPACITY_IN_POUNDS Imputation (KNN with features)')
62 ax.legend()
63 plt.show()
```

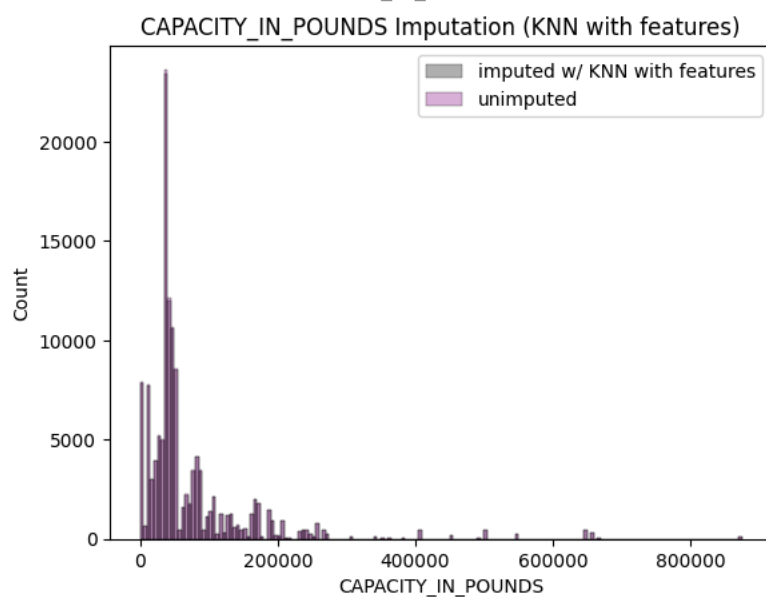
RMSE (Median - CAPACITY\_IN\_POUNDS): 92646.40



RMSE (KNN - CAPACITY\_IN\_POUNDS): 92646.40



RMSE (KNN with features - CAPACITY\_IN\_POUNDS): 44381.46



MANUFACTURE\_YEAR:

We can impute MANUFACTURE\_YEAR because it is a numeric, discrete variable, the missing values are relatively rare and likely missing at random, and there's a clear association with other variables like MODEL or AIRLINE\_ID.

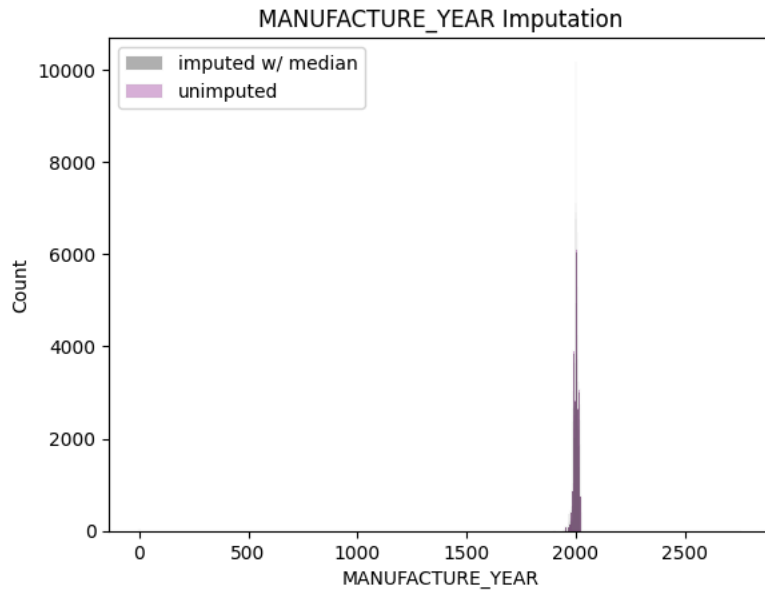
I chose to use Median because it works well when data is clustered and I chose PMM because it guarantees realistic values and avoids artificial midpoints.

```

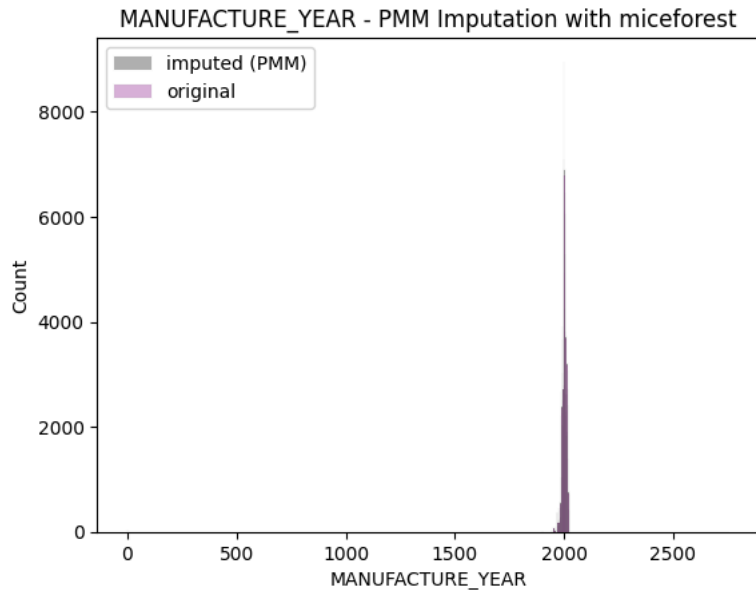
1 # random missing
2 df_temp = df[['MANUFACTURE_YEAR']].dropna().copy()
3 df_temp.loc[df_temp.sample(frac=0.01, random_state=0).index, 'MANUFACTURE_YEAR'] = np.nan
4
5 true_vals_year = df[['MANUFACTURE_YEAR']].dropna().loc[df_temp[df_temp['MANUFACTURE_YEAR'].isna()].index]
6
7 # Median
8 median_imputer_year = SimpleImputer(strategy='median')
9 df_temp['MANUFACTURE_YEAR'] = median_imputer_year.fit_transform(df_temp[['MANUFACTURE_YEAR']])
10
11 # RMSE
12 rmse_year = np.sqrt(mean_squared_error(true_vals_year, df_temp.loc[true_vals_year.index]))
13 print(f"RMSE (MANUFACTURE_YEAR): {rmse_year:.2f}")
14
15 # Plot
16 fig, ax = plt.subplots(1, 1)
17 sns.histplot(df_temp['MANUFACTURE_YEAR'], binwidth=1, alpha=0.3, color='k', label='imputed w/ median')
18 sns.histplot(df['MANUFACTURE_YEAR'], binwidth=1, alpha=0.3, color='purple', label='unimputed')
19 ax.set_title('MANUFACTURE_YEAR Imputation')
20 ax.legend()
21 plt.show()
22
23 import miceforest as mf
24
25 # PMM
26 df_pmm = df[['MANUFACTURE_YEAR', 'CAPACITY_IN_POUNDS', 'NUMBER_OF_SEATS']].dropna().reset_index(drop=True)
27
28 missing_idx = df_pmm.sample(frac=0.01, random_state=0).index
29 true_vals = df_pmm.loc[missing_idx, 'MANUFACTURE_YEAR'].copy()
30 df_pmm.loc[missing_idx, 'MANUFACTURE_YEAR'] = np.nan
31
32 num_datasets = 4
33 kernel = mf.ImputationKernel(
34     data=df_pmm,
35     num_datasets=num_datasets,
36     save_all_iterations_data=False,
37     random_state=1
38 )
39
40 kernel.mice(1)
41
42 df_completed = kernel.complete_data(dataset=0)
43
44 # RMSE
45 rmse_pmm = np.sqrt(mean_squared_error(true_vals, df_completed.loc[missing_idx, 'MANUFACTURE_YEAR']))
46 print(f"RMSE (PMM - MANUFACTURE_YEAR): {rmse_pmm:.2f}")
47
48 # Plot
49 fig, ax = plt.subplots(1, 1)
50 sns.histplot(df_completed['MANUFACTURE_YEAR'], binwidth=1, alpha=0.3, color='k', label='imputed (PMM)')
51 sns.histplot(df_pmm['MANUFACTURE_YEAR'].dropna(), binwidth=1, alpha=0.3, color='purple', label='original')
52 ax.set_title('MANUFACTURE_YEAR - PMM Imputation with miceforest')
53 ax.legend()
54 plt.show()
55
56

```

RMSE (MANUFACTURE\_YEAR): 10.16



RMSE (PMM - MANUFACTURE\_YEAR): 7.06



NUMBER\_OF\_SEATS:

We can impute NUMBER\_OF\_SEATS because it's a discrete, positive integer tied to the MODEL of the aircraft and there's low variation within the model.

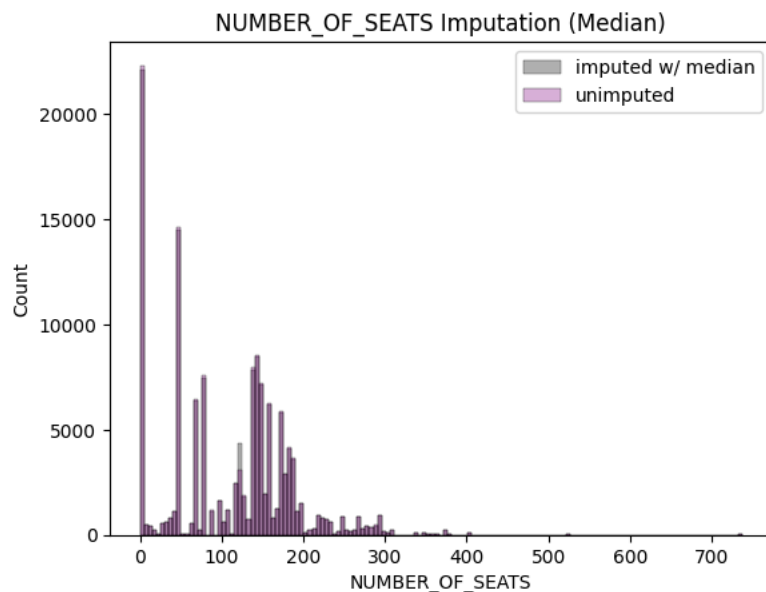
I chose to use Median because it's simple and most aircrafts of the same model have the same seat amount. Same as with CAPACITY\_IN\_POUNDS I found that KNN alone was similar to Median but KNN with features improved accuracy by taking into account multiple features.

```
1 # random missing
2 df_temp = df[['NUMBER_OF_SEATS']].dropna().copy()
3 df_temp.loc[df_temp.sample(frac=0.01, random_state=0).index, 'NUMBER_OF_SEATS'] = np.nan
4
5 true_vals_seats = df[['NUMBER_OF_SEATS']].dropna().loc[df_temp[df_temp['NUMBER_OF_SEATS'].isna()].index]
6
7 # Median
8 median_imputer_seats = SimpleImputer(strategy='median')
9 df_temp['NUMBER_OF_SEATS'] = median_imputer_seats.fit_transform(df_temp[['NUMBER_OF_SEATS']])
10
11 # RMSE
12 rmse_seats = np.sqrt(mean_squared_error(true_vals_seats, df_temp.loc[true_vals_seats.index]))
13 print(f"RMSE (Median - NUMBER_OF_SEATS): {rmse_seats:.2f}")
14
15 # Plot
```

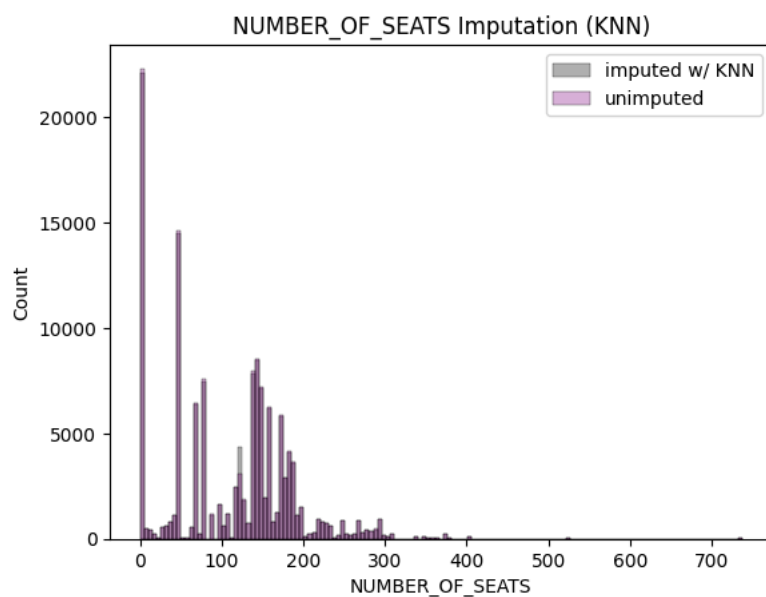
```
16 fig, ax = plt.subplots(1, 1)
17 sns.histplot(df_temp['NUMBER_OF_SEATS'], binwidth=5, alpha=0.3, color='k', label='imputed w/ median')
18 sns.histplot(df['NUMBER_OF_SEATS'], binwidth=5, alpha=0.3, color='purple', label='unimputed')
19 ax.set_title('NUMBER_OF_SEATS Imputation (Median)')
20 ax.legend()
21 plt.show()
22
23 # KNN
24 knn_imputer_seats = KNNImputer(n_neighbors=5)
25 df_temp['NUMBER_OF_SEATS'] = knn_imputer_seats.fit_transform(df_temp[['NUMBER_OF_SEATS']])
26
27 # RMSE
28 rmse_seats_knn = np.sqrt(mean_squared_error(true_vals_seats, df_temp.loc[true_vals_seats.index]))
29 print(f"RMSE (KNN - NUMBER_OF_SEATS): {rmse_seats_knn:.2f}")
30
31 # Plot
32 fig, ax = plt.subplots(1, 1)
33 sns.histplot(df_temp['NUMBER_OF_SEATS'], binwidth=5, alpha=0.3, color='k', label='imputed w/ KNN')
34 sns.histplot(df['NUMBER_OF_SEATS'], binwidth=5, alpha=0.3, color='purple', label='unimputed')
35 ax.set_title('NUMBER_OF_SEATS Imputation (KNN)')
36 ax.legend()
37 plt.show()
38
39 # KNN with features
40 knn_df = df[['NUMBER_OF_SEATS', 'CAPACITY_IN_POUNDS', 'MANUFACTURE_YEAR']].dropna().reset_index(drop=True)
41
42 missing_idx = knn_df.sample(frac=0.01, random_state=0).index
43 true_vals = knn_df.loc[missing_idx, 'NUMBER_OF_SEATS'].copy()
44 knn_df.loc[missing_idx, 'NUMBER_OF_SEATS'] = np.nan
45
46 knn_imputer = KNNImputer(n_neighbors=5)
47 knn_imputed = pd.DataFrame(knn_imputer.fit_transform(knn_df), columns=knn_df.columns)
48
49 # RMSE
50 rmse_seats_knn_features = np.sqrt(mean_squared_error(true_vals, knn_imputed.loc[missing_idx, 'NUMBER_OF_SEATS']))
51 print(f"RMSE (KNN with features - NUMBER_OF_SEATS): {rmse_seats_knn_features:.2f}")
52
53 # Plot
54 fig, ax = plt.subplots(1, 1)
55 sns.histplot(knn_imputed['NUMBER_OF_SEATS'], binwidth=5, alpha=0.3, color='k', label='imputed w/ KNN with features')
56 sns.histplot(knn_df['NUMBER_OF_SEATS'], binwidth=5, alpha=0.3, color='purple', label='original')
57 ax.set_title('NUMBER_OF_SEATS Imputation (KNN with features)')
58 ax.legend()
59 plt.show()
```



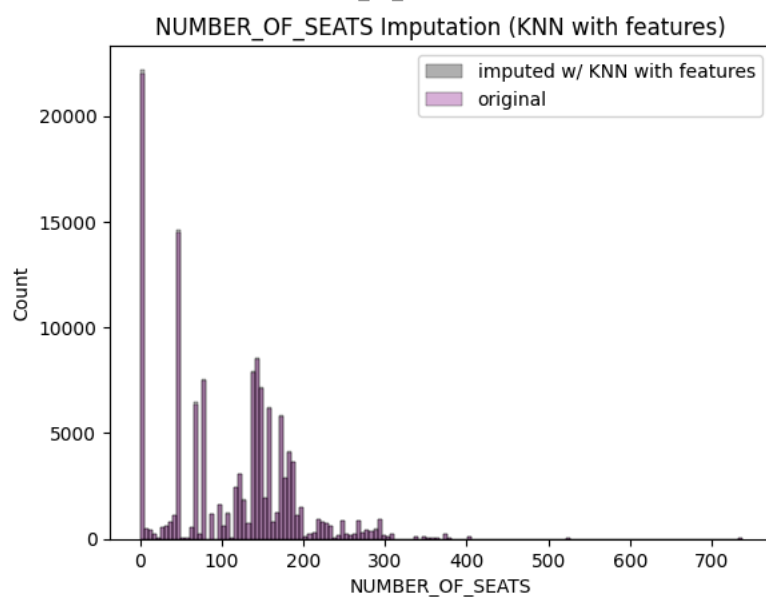
RMSE (Median - NUMBER\_OF\_SEATS): 77.11



RMSE (KNN - NUMBER\_OF\_SEATS): 77.11



RMSE (KNN with features - NUMBER\_OF\_SEATS): 16.60



```

1 # Adding numerical imputed columns to the df
2 df_extended = df.copy()
3
4 df_extended.loc[df_completed.index, 'MANUFACTURE_YEAR_PMM'] = df_completed['MANUFACTURE_YEAR'].round()
5 df_extended.loc[knn_df.index, 'NUMBER_OF_SEATS_KNN_FEATURES'] = knn_imputed['NUMBER_OF_SEATS'].round()
6 df_extended.loc[knn_df.index, 'CAPACITY_IN_POUNDS_KNN_FEATURES'] = knn_imputed['CAPACITY_IN_POUNDS'].round()
7

```

## ✓ Task 2

Inspect the columns MANUFACTURER, MODEL, AIRCRAFT\_STATUS, and OPERATING\_STATUS and decide if transformation or standardization of data are required.

MANUFACTURER column

```

1 # convert to upper and get rid of spaces
2 df_extended['MANUFACTURER_CLEAN'] = df_extended['MANUFACTURER'].str.upper().str.strip()
3
4 manufacturer_map = {
5     "BOEING": "BOEING",
6     "THEBOEINGCO": "BOEING",
7     "BOEING COMPANY": "BOEING",
8     "BOEINGCO": "BOEING",
9     "AIRBUS INDUSTRIE": "AIRBUS",
10    "AIRBUS": "AIRBUS",
11    "BOMBARDIER INC": "BOMBARDIER",
12    "EMBRAER S.A.": "EMBRAER",
13    "EMBRAER": "EMBRAER",
14    "GULFSTREAM AEROSPACE": "GULFSTREAM",
15    "CESSNA AIRCRAFT COMPANY": "CESSNA",
16    "CESSNA": "CESSNA"
17 }
18
19 df_extended['MANUFACTURER_STD'] = df_extended['MANUFACTURER_CLEAN'].replace(manufacturer_map)
20
21 print(df_extended['MANUFACTURER_STD'].value_counts().head(20))

```

```

MANUFACTURER_STD
BOEING                44191
EMBRAER               15554
AIRBUS                13440
BOMBARDIER            11834
AIRBUSINDUSTRIES      7053
BOEINGCOMPANY         6767
CESSNA                4514
MCDONNELLDUGLAS       4306
MCDONNELL-DOUGLAS     4159
THEBOEINGCOMPANY      3975
CANADAIR              3861
AIRBUSINDUSTRIE       2666
ATR                   1181
DOUGLAS               1137
GE                    1110
DEHAVILLAND           1084
MCDONNELLDUGLAS       736
BOMBARDIERAEROSPACE   649
BOEINGCO.             584
GULFSTREAM            441
Name: count, dtype: int64

```

MODEL column

```

1 # convert to upper and get rid of spaces
2 df_extended['MODEL_CLEAN'] = df_extended['MODEL'].str.upper().str.strip()
3
4 print(df_extended['MODEL_CLEAN'].value_counts().head(30))
5
6 model_map = {
7     "737-800": "BOEING 737-800",
8     "737": "BOEING 737",
9     "B737": "BOEING 737",
10    "B737-800": "BOEING 737-800",
11    "A320": "AIRBUS A320",

```

```

12 "A320-200": "AIRBUS A320",
13 "BOMBARDIER CRJ200": "CRJ200",
14 "CRJ-200": "CRJ200",
15 "EMBRAER 170": "EMBRAER E170",
16 "EMB-170": "EMBRAER E170"
17 }
18
19 df_extended['MODEL_STD'] = df_extended['MODEL_CLEAN'].replace(model_map)
20
21 print(df_extended['MODEL_STD'].value_counts().head(10))
22

```

```

MODEL_CLEAN
EMB-145                2976
B-737-7H4              2470
B737-823               2370
A320-232               2333
A321-231               2259
737-700PASSENGERONLY  2027
C-208B                 1872
B757-2                 1775
CRJ-2/4                1761
B737-800PAX            1621
MD-80                  1610
A320-1/2               1466
ERJ-170-200LR          1379
B737-7/L               1349
757-200                1345
CRJ200-2B19            1342
A319                   1267
B-737-8H4              1256
CRJ-200                1148
ERJ-175                1132
SUPER80PASSENGER       1108
MD-11                  1107
757-24APF              1039
B737-3                 1036
MD-88-PSGR             1028
C-208B/3               1017
757-232-PSGR           988
CRJ-900LR-PSGR          976
B737-823PASSENGER       956
B737-8                  951
Name: count, dtype: int64
MODEL_STD
EMB-145                2976
B-737-7H4              2470
B737-823               2370
A320-232               2333
A321-231               2259
737-700PASSENGERONLY  2027
CRJ200                 1929
C-208B                 1872
B757-2                 1775
CRJ-2/4                1761
Name: count, dtype: int64

```

AIRCRAFT\_STATUS column

```

1 # convert to lower and get rid of spaces
2 df_extended['AIRCRAFT_STATUS_CLEAN'] = df_extended['AIRCRAFT_STATUS'].str.lower().str.strip()
3
4 print(df_extended['AIRCRAFT_STATUS_CLEAN'].value_counts())
5
6 aircraft_status_map = {
7     "active": "active",
8     "inactive": "inactive",
9     "retired": "retired",
10    "destroyed": "retired",
11    "decommissioned": "retired",
12    "in service": "active",
13    "not in service": "inactive"
14 }
15
16 df_extended['AIRCRAFT_STATUS_STD'] = df_extended['AIRCRAFT_STATUS_CLEAN'].replace(aircraft_status_map)
17
18 print(df_extended['AIRCRAFT_STATUS_STD'].value_counts())

```

```

AIRCRAFT_STATUS_CLEAN
0      79506

```

```

b    43551
a     9134
l      122
Name: count, dtype: int64
AIRCRAFT_STATUS_STD
o    79506
b    43551
a     9134
l      122
Name: count, dtype: int64

```

OPERATING\_STATUS column

```

1 # convert to lower and get rid of spaces
2 df_extended['OPERATING_STATUS_CLEAN'] = df_extended['OPERATING_STATUS'].str.lower().str.strip()
3
4 print(df_extended['OPERATING_STATUS_CLEAN'].value_counts())
5
6 operating_status_map = {
7     "operating": "operating",
8     "not operating": "not operating",
9     "inactive": "not operating",
10    "active": "operating",
11    "temporarily inactive": "not operating",
12    "unknown": "unknown"
13 }
14
15 df_extended['OPERATING_STATUS_STD'] = df_extended['OPERATING_STATUS_CLEAN'].replace(operating_status_map)
16
17 print(df_extended['OPERATING_STATUS_STD'].value_counts())
18
19

```

```

↗ OPERATING_STATUS_CLEAN
y    126648
n     5664
      1
Name: count, dtype: int64
OPERATING_STATUS_STD
y    126648
n     5664
      1
Name: count, dtype: int64

```

### Task 3

Remove data rows with missing values.

```

1 cleaned_aircraft_inventory = df_extended.dropna()
2 print(f"Remaining rows after dropping missing values: {len(cleaned_aircraft_inventory)}")

```

```

↗ Remaining rows after dropping missing values: 101096

```

### Task 4

Transformation and derivative variables with NUMBER\_OF\_SEATS and CAPACITY\_IN\_POUNDS

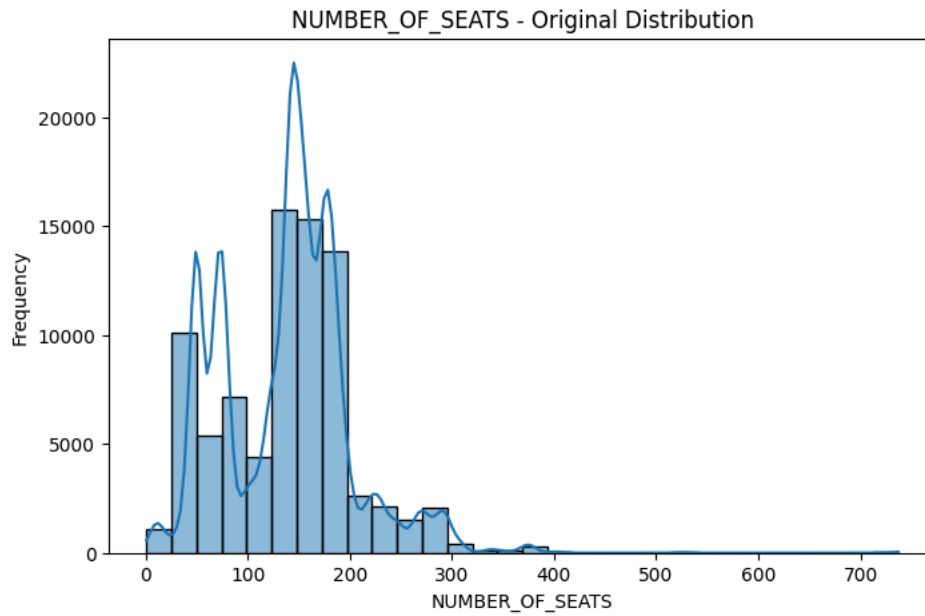
```

1 from scipy.stats import boxcox
2
3 filtered_data = cleaned_aircraft_inventory[(cleaned_aircraft_inventory['NUMBER_OF_SEATS'] > 0) &
4                                             (cleaned_aircraft_inventory['CAPACITY_IN_POUNDS'] > 0)].copy()
5
6 # check skewness and plot original histograms
7 for col in ['NUMBER_OF_SEATS', 'CAPACITY_IN_POUNDS']:
8     skewness = filtered_data[col].skew()
9     print(f"{col} skewness before Box-Cox: {skewness:.2f}")
10
11 plt.figure(figsize=(8, 5))
12 sns.histplot(filtered_data[col], bins=30, kde=True)
13 plt.title(f'{col} - Original Distribution')
14 plt.xlabel(col)

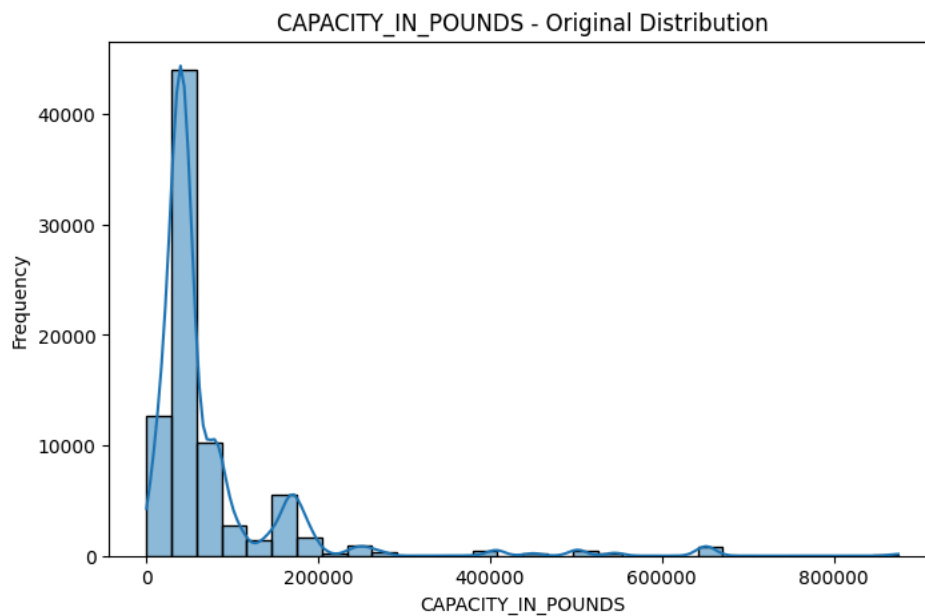
```

```
15     plt.ylabel('Frequency')
16     plt.show()
17
18 # apply Box-Cox transformation
19 filtered_data['NUMBER_OF_SEATS_BOXCOX'], _ = boxcox(filtered_data['NUMBER_OF_SEATS'])
20 filtered_data['CAPACITY_IN_POUNDS_BOXCOX'], _ = boxcox(filtered_data['CAPACITY_IN_POUNDS'])
21
22 # check skewness and plot new histograms
23 for col in ['NUMBER_OF_SEATS_BOXCOX', 'CAPACITY_IN_POUNDS_BOXCOX']:
24     skewness = filtered_data[col].skew()
25     print(f"{col} skewness after Box-Cox: {skewness:.2f}")
26
27     plt.figure(figsize=(8, 5))
28     sns.histplot(filtered_data[col], bins=30, kde=True)
29     plt.title(f'{col} - After Box-Cox Transformation')
30     plt.xlabel(col)
31     plt.ylabel('Frequency')
32     plt.show()
33
```

NUMBER\_OF\_SEATS skewness before Box-Cox: 0.78



CAPACITY\_IN\_POUNDS skewness before Box-Cox: 4.18



NUMBER\_OF\_SEATS\_BOXCOX skewness after Box-Cox: 0.01

