

## ▼ NOAA Dataset Cleaning

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force\_remount=True).

```
1 import pandas as pd
2 import numpy as np
3 from IPython.display import display
4 import matplotlib.pyplot as plt
5 import seaborn as sns
```

```
1 noaa_weather_raw = pd.read_csv("/content/drive/SharedDrives/Data Science 303 Group Project/csv/noaa_weather/noaa_CA_1992_2016_weather_2781174.csv")
2 print(noaa_weather_raw.info())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8761 entries, 0 to 8760
Data columns (total 29 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   STATION      8761 non-null   object
1   NAME         8761 non-null   object
2   LATITUDE     8761 non-null   float64
3   LONGITUDE    8761 non-null   float64
4   ELEVATION    8761 non-null   float64
5   DATE         8761 non-null   object
6   CDSO         7409 non-null   float64
7   CLDD         8683 non-null   float64
8   DP01         266 non-null    float64
9   DP10         266 non-null    float64
10  DSND         258 non-null    float64
11  DSNW         247 non-null    float64
12  DT00         8724 non-null   float64
13  DT32         8724 non-null   float64
14  DX32         8738 non-null   float64
15  DX70         8738 non-null   float64
16  DX90         8738 non-null   float64
17  EMNT         8724 non-null   float64
18  EMSD         258 non-null    float64
19  EMSN         247 non-null    float64
20  EMXP         266 non-null    float64
21  EMXT         8738 non-null   float64
22  HDSO         7594 non-null   float64
23  HTDD         8683 non-null   float64
24  PRCP         266 non-null    float64
25  SNOW         247 non-null    float64
26  TAVG         8705 non-null   float64
27  TMAX         8738 non-null   float64
28  TMIN         8724 non-null   float64
dtypes: float64(26), object(3)
memory usage: 1.9+ MB
None
```

```
1 NOAA_DECODE = {
2     "CDSO": "NUM_COOLING_DEGREE_DAY_CUMULATIVE",
3     "CLDD": "NUM_COOLING_DEGREE_DAY",
4     "DP01": "NUM_DAYS_WITH_0_01_INCH_PRECIPITATION",
5     "DP10": "NUM_DAYS_WITH_0_1_INCH_PRECIPITATION",
6     "DSND": "NUM_DAYS_SNOW_DEPTH_1_INCH",
7     "DSNW": "NUM_DAYS_SNOW_FALL_1_INCH",
8     "DT00": "NUM_DAYS_WITH_MIN_TEMP_BELOW_0_FAHRENHEIT",
9     "DT32": "NUM_DAYS_WITH_MIN_TEMP_BELOW_32_FAHRENHEIT",
10    "DX70": "NUM_DAYS_WITH_MAX_TEMP_ABOVE_70_FAHRENHEIT",
11    "DX90": "NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT",
12    "EMNT": "EXTREME_MINIMUM_TEMPERATURE_FOR_MONTH",
13    "EMSD": "HIGHEST_DAILY_SNOW_DEPTH",
14    "EMSN": "HIGHEST_DAILY_SNOW_FALL",
15    "EMXP": "HIGHEST_DAILY_PRECIPITATION",
16    "EMXT": "EXTREME_MAXIMUM_TEMPERATURE_MONTH",
17    "HDSO": "HEATING_DEGREE_DAYS_TO_DATE",
18    "HTDD": "NUM_DAYS WHERE AVG_TEMP_BELOW_65_FAHRENHEIT",
19    "PRCP": "TOTAL_MONTHLY_RAINFALL",
20    "SNOW": "TOTAL_MONTHLY_SNOWFALL",
21    "TAVG": "TEMPERATURE_AVERAGE",
22    "TMAX": "TEMPERATURE_MAX",
23    "TMIN": "TEMPERATURE_MIN"
```

```
24 }
25
26
27 # DECODE AND ORGANIZE NOAA_DATA
28 noaa_weather = noaa_weather_raw[["STATION", "LATITUDE", "LONGITUDE", "ELEVATION"]]
29
30 # CONVERT NAME OF STATION TO UNDERScores
31 noaa_weather["STATION_NAME"] = noaa_weather_raw.NAME.map(lambda raw_name: raw_name[: (raw_name.index(" CALI"))].replace(" ", "_"))
32
33 # CONVERT DATE STRING TO DATE TIME YEAR AND MONTH
34 dates = pd.to_datetime(noaa_weather_raw["DATE"])
35 year = dates.dt.year
36 month = dates.dt.month
37
38 noaa_weather["YEAR"] = year
39 noaa_weather["MONTH"] = month
40
41
42 # CONVERT COLUMN NAMES TO ENGLISH
43 for (k, v) in NOAA_DECODE.items():
44     length = len(noaa_weather_raw[k])
45     num_missing = noaa_weather_raw[k].isnull().sum()
46     print(f"Feature '{v}' is missing {(num_missing / length)*100:1.1f}%")
47
48 # Exclude features that have > 90% of values missing
49 if not (num_missing / length >= .9):
50     noaa_weather[v] = noaa_weather_raw[k]
51
52 display(noaa_weather)
53 noaa_weather.info()
```

```
Feature 'NUM_COOLING_DEGREE_DAY_CUMULATIVE' is missing 15.4%
Feature 'NUM_COOLING_DEGREE_DAY' is missing 0.9%
Feature 'NUM_DAYS_WITH_0_01_INCH_PRECIPITATION' is missing 97.0%
Feature 'NUM_DAYS_WITH_0_1_INCH_PRECIPITATION' is missing 97.0%
Feature 'NUM_DAYS_SNOW_DEPTH_1_INCH' is missing 97.1%
Feature 'NUM_DAYS_SNOW_FALL_1_INCH' is missing 97.2%
Feature 'NUM_DAYS_WITH_MIN_TEMP_BELOW_0_FAHRENHEIT' is missing 0.4%
Feature 'NUM_DAYS_WITH_MIN_TEMP_BELOW_32_FAHRENHEIT' is missing 0.4%
Feature 'NUM_DAYS_WITH_MAX_TEMP_ABOVE_70_FAHRENHEIT' is missing 0.3%
Feature 'NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT' is missing 0.3%
Feature 'EXTREME_MINIMUM_TEMPERATURE_FOR_MONTH' is missing 0.4%
Feature 'HIGHEST_DAILY_SNOW_DEPTH' is missing 97.1%
Feature 'HIGHEST_DAILY_SNOW_FALL' is missing 97.2%
Feature 'HIGHEST_DAILY_PRECIPITATION' is missing 97.0%
Feature 'EXTREME_MAXIMUM_TEMPERATURE_MONTH' is missing 0.3%
Feature 'HEATING_DEGREE_DAYS_TO_DATE' is missing 13.3%
Feature 'NUM_DAYS_WHERE_AVG_TEMP_BELOW_65_FAHRENHEIT' is missing 0.9%
Feature 'TOTAL_MONTHLY_RAINFALL' is missing 97.0%
Feature 'TOTAL_MONTHLY_SNOWFALL' is missing 97.2%
Feature 'TEMPERATURE_AVERAGE' is missing 0.6%
Feature 'TEMPERATURE_MAX' is missing 0.3%
Feature 'TEMPERATURE_MIN' is missing 0.4%
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:31: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

```
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:38: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

```
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:39: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

```
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:50: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

```
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
```

	STATION	LATITUDE	LONGITUDE	ELEVATION	STATION_NAME	YEAR	MONTH	NUM_COOLING_DEGREE_DAY_CUMULATIVE	NUM_COOLING_DEGREE_DAY	NUM_DAYS_WITH_MIN_TEMP_BELOW_0
0	USR0000CCOH	39.8717	-121.7689	528.2	COHASSET	1992	1	0.0		0.0

Impute Values

1	2	1	1	1	1	1	1	1	1	1
---	---	---	---	---	---	---	---	---	---	---

```
1 # First, we want to drop stations that don't have enough data over the given time period
2 max_year = 2016
3 min_year = noaa_weather['YEAR'].min()
4 station_name_list = list(noaa_weather["STATION_NAME"].unique())
5 total_station_month_year_pairs_expected = (max_year - min_year + 1) * 12
6
7 # Print out the percentage of missing station, month, year pairs:
8 station_drop_list = []
9 for station in station_name_list:
10     num_present = (noaa_weather["STATION_NAME"] == station).sum()
11     # print(f'station: {station}, { num_present / total_station_month_year_pairs_expected}')
12     if num_present / total_station_month_year_pairs_expected < .75:
13         station_drop_list.append(station)
14
15 # For any station with less than 75 percent of month-date pairs missing, drop them
16 for station in station_drop_list:
17     print(station)
18     noaa_weather.drop(noaa_weather[noaa_weather['STATION_NAME'] == station].index, inplace = True)
```

```
/usr/local/lib/python3.7/dist-packages/pandas/core/frame.py:4174: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame
```

```
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
errors=errors,
MALLORY RIDGE
BRADLEY
CASHMAN
BIG ROCK
MAD RIVER
EEL RIVER CAMP
JOHNSONDALE
MODOC_NWR
AMMO_DUMP
TEMESCAL_LPF
LYTLE_CREEK
MAPLE_CREEK
```

FAWNSKIN  
CLEAR\_CREEK  
BEVERLY\_HILLS  
MILO  
KERNVILLE  
TRIMMER  
CARPENTER\_RIDGE  
OJAI  
CRANSTON

```
1 def count_missing_data(df):
2     # Count the number of missing values by feature:
3     df_na = df.isna().sum()
4     df_na = df_na.drop(df_na[df_na == 0].index).sort_values(ascending=False)
5     df_na = (df_na / len(df)) * 100
6     missing_data = pd.DataFrame({"Missing Ratio" : df_na})
7     display(missing_data)
8
9 count_missing_data(noaa_weather)
```

	Missing Ratio
NUM_COOLING_DEGREE_DAY_CUMULATIVE	15.925234
HEATING_DEGREE_DAYS_TO_DATE	14.074766
NUM_DAYS_WHERE_AVG_TEMP_BELOW_65_FAHRENHEIT	0.934579
NUM_COOLING_DEGREE_DAY	0.934579
TEMPERATURE_AVERAGE	0.598131
TEMPERATURE_MAX	0.392523
EXTREME_MAXIMUM_TEMPERATURE_MONTH	0.392523
NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT	0.392523
NUM_DAYS_WITH_MAX_TEMP_ABOVE_70_FAHRENHEIT	0.392523
TEMPERATURE_MIN	0.280374
EXTREME_MINIMUM_TEMPERATURE_FOR_MONTH	0.280374
NUM_DAYS_WITH_MIN_TEMP_BELOW_32_FAHRENHEIT	0.280374

```
1 # We will next impute values. We want to impute only from other samples from the same weather station. We will group by station AND month.
2 # We can expect that year to year, a value should be about the same for each month
3 stations = noaa_weather.STATION_NAME.unique()
4 features_by_station_and_month = noaa_weather.groupby([noaa_weather.STATION_NAME, noaa_weather.MONTH])
5 for feature in noaa_weather.select_dtypes(include="number").columns:
6     # Skip imputing the year, month, latitude, or longitude:
7     if not feature == "YEAR" and not feature == "MONTH" and not feature == "LATITUDE" and not feature == "LONGITUDE":
8         print(f"Imputing feature by station and month { feature }")
9         noaa_weather[feature] = features_by_station_and_month[feature].transform(lambda group: group.fillna(group.mean()))
10
11 count_missing_data(noaa_weather)
```

Imputing feature by station and month ELEVATION  
Imputing feature by station and month NUM\_COOLING\_DEGREE\_DAY\_CUMULATIVE  
/usr/local/lib/python3.7/dist-packages/ipykernel\_launcher.py:9: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
if __name__ == '__main__':
    Imputing feature by station and month NUM_COOLING_DEGREE_DAY
    Imputing feature by station and month NUM_DAYS_WITH_MIN_TEMP_BELOW_0_FAHRENHEIT
    Imputing feature by station and month NUM_DAYS_WITH_MIN_TEMP_BELOW_32_FAHRENHEIT
    Imputing feature by station and month NUM_DAYS_WITH_MAX_TEMP_ABOVE_70_FAHRENHEIT
    Imputing feature by station and month NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT
    Imputing feature by station and month EXTREME_MINIMUM_TEMPERATURE_FOR_MONTH
    Imputing feature by station and month EXTREME_MAXIMUM_TEMPERATURE_MONTH
    Imputing feature by station and month HEATING_DEGREE_DAYS_TO_DATE
    Imputing feature by station and month NUM_DAYS_WHERE_AVG_TEMP_BELOW_65_FAHRENHEIT
    Imputing feature by station and month TEMPERATURE_AVERAGE
    Imputing feature by station and month TEMPERATURE_MAX
    Imputing feature by station and month TEMPERATURE_MIN

    Missing Ratio
```

```
1 # Check to make sure all columns have correct type.... Looks good
2 noaa_weather.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 5350 entries, 0 to 8448
Data columns (total 20 columns):
#   Column                                     Non-Null Count  Dtype
---  -
0   STATION                                     5350 non-null   object
1   LATITUDE                                   5350 non-null   float64
2   LONGITUDE                                   5350 non-null   float64
3   ELEVATION                                   5350 non-null   float64
4   STATION_NAME                               5350 non-null   object
5   YEAR                                        5350 non-null   int64
6   MONTH                                       5350 non-null   int64
7   NUM_COOLING_DEGREE_DAY_CUMULATIVE         5350 non-null   float64
8   NUM_COOLING_DEGREE_DAY                     5350 non-null   float64
9   NUM_DAYS_WITH_MIN_TEMP_BELOW_0_FAHRENHEIT  5350 non-null   float64
10  NUM_DAYS_WITH_MIN_TEMP_BELOW_32_FAHRENHEIT 5350 non-null   float64
11  NUM_DAYS_WITH_MAX_TEMP_ABOVE_70_FAHRENHEIT 5350 non-null   float64
12  NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT 5350 non-null   float64
13  EXTREME_MINIMUM_TEMPERATURE_FOR_MONTH       5350 non-null   float64
14  EXTREME_MAXIMUM_TEMPERATURE_MONTH           5350 non-null   float64
15  HEATING_DEGREE_DAYS_TO_DATE                 5350 non-null   float64
16  NUM_DAYS_WHERE_AVG_TEMP_BELOW_65_FAHRENHEIT 5350 non-null   float64
17  TEMPERATURE_AVERAGE                       5350 non-null   float64
18  TEMPERATURE_MAX                             5350 non-null   float64
19  TEMPERATURE_MIN                             5350 non-null   float64
dtypes: float64(16), int64(2), object(2)
memory usage: 877.7+ KB
```

```
1 # TODO: CHECK IF THIS GIVES GOOD RESULTS. HAVING TO IMPUTE A LOT!!
2
3
4 # Check which station, month, year pairs are missing
5 import itertools
6
7 noaa_stations = list(noaa_weather["STATION_NAME"].unique())
8 years = [i for i in range(1992, 2016)]
9 months = [i for i in range(1, 13)]
10
11 station_month_year_pairs = list(itertools.product(noaa_stations, months, years))
12 # print(station_month_year_pairs)
13 pairs_missing = 0
14 pairs_not_missing = 0
15
16 all_missing_triplets = []
17 for triple in station_month_year_pairs:
18     if not ((noaa_weather["STATION_NAME"] == triple[0]) & (noaa_weather["YEAR"] == triple[2]) & (noaa_weather["MONTH"] == triple[1])).any():
19         # print(f"Missing: {triple[0]}:{triple[1]}:{triple[2]}")
20         all_missing_triplets.append((triple))
21         pairs_missing += 1
22     else:
23         pairs_not_missing += 1
24
25 print(f"Total Station_Month_Year triplets missing: {pairs_missing}, or {pairs_missing / (pairs_missing + pairs_not_missing) * 100}%")
```

Total Station\_Month\_Year triplets missing: 429, or 7.447916666666667%

```
1 # With the given list of missing station, month, year pairs, impute the values.
2 print("Will now impute the prior values")
3 def impute_missing_month_row_year_pairs(true_records, imputed_records, station, month, year):
4     """
5     Given an incomplete dataset, and a missing (station, year, month) triplet,
6     impute a row using the available data.
7
8     Hierarchy to generate / impute a value:
9     1. Reuse last recorded station month year pair. For example, if ("A", 2000, 7) is missing, fill it in with ("A", 1999, 7)
10    2. If no years prior to the missing year exist, use the least recent recording from the same station
11    """
12    matching = true_records[((true_records["STATION_NAME"] == station) & (true_records["MONTH"] == month))]
13    months_prior_to_year = matching[matching["YEAR"] < year]
14
15    if months_prior_to_year.shape[0]:
16        idx = months_prior_to_year["YEAR"].argmax()
17        best_match = months_prior_to_year.iloc[idx].copy()
18        print(f"PAST MATCH: { station }, { month }, { year } -> {best_match.STATION_NAME}, {best_match.MONTH}, {best_match.YEAR}")
19    else:
20        # Get the least recent year if no prior years exist
21        months_after_year = matching[matching["YEAR"] > year]
```

```

42     idx = montns_after_year[`YEAR`.j.argmax()]
23     best_match = months_after_year.iloc[idx].copy()
24     print(f"FUTURE MATCH: { station }, { month }, { year } -> {best_match.STATION_NAME}, {best_match.MONTH}, {best_match.YEAR}" )
25
26     best_match[ "YEAR" ] = year
27     best_match[ "MONTH" ] = month
28
29     imputed_records.append(best_match)
30
31 imputed_records = []
32 for missing_record in all_missing_triplets:
33     impute_missing_month_row_year_pairs(noaa_weather, imputed_records, missing_record[0], missing_record[1], missing_record[2])
34
35 noaa_weather = noaa_weather.append(imputed_records)
36 # get_matching_station(noaa_weather, "HUNTER MOUNTAIN", 12, 2012)

```

```

Will now impute the prior values
PAST MATCH: COHASSET, 1, 1995 -> COHASSET, 1, 1994
PAST MATCH: COHASSET, 2, 1995 -> COHASSET, 2, 1994
PAST MATCH: COHASSET, 5, 2012 -> COHASSET, 5, 2011
PAST MATCH: COHASSET, 6, 2012 -> COHASSET, 6, 2011
PAST MATCH: COHASSET, 12, 1994 -> COHASSET, 12, 1993
PAST MATCH: LADDER BUTTE, 1, 1993 -> LADDER BUTTE, 1, 1992
PAST MATCH: LADDER BUTTE, 1, 1995 -> LADDER BUTTE, 1, 1994
PAST MATCH: LADDER BUTTE, 1, 2000 -> LADDER BUTTE, 1, 1999
PAST MATCH: LADDER BUTTE, 2, 1993 -> LADDER BUTTE, 2, 1992
PAST MATCH: LADDER BUTTE, 2, 1995 -> LADDER BUTTE, 2, 1994
PAST MATCH: LADDER BUTTE, 3, 1993 -> LADDER BUTTE, 3, 1992
PAST MATCH: LADDER BUTTE, 3, 1995 -> LADDER BUTTE, 3, 1994
PAST MATCH: LADDER BUTTE, 4, 1993 -> LADDER BUTTE, 4, 1992
PAST MATCH: LADDER BUTTE, 4, 1995 -> LADDER BUTTE, 4, 1994
PAST MATCH: LADDER BUTTE, 11, 1995 -> LADDER BUTTE, 11, 1994
PAST MATCH: LADDER BUTTE, 11, 2000 -> LADDER BUTTE, 11, 1999
FUTURE MATCH: LADDER BUTTE, 12, 1992 -> LADDER BUTTE, 12, 1993
PAST MATCH: JUANITA LAKE, 1, 1994 -> JUANITA LAKE, 1, 1993
PAST MATCH: JUANITA LAKE, 1, 2000 -> JUANITA LAKE, 1, 1999
PAST MATCH: JUANITA LAKE, 1, 2001 -> JUANITA LAKE, 1, 1999
PAST MATCH: JUANITA LAKE, 1, 2002 -> JUANITA LAKE, 1, 1999
PAST MATCH: JUANITA LAKE, 1, 2003 -> JUANITA LAKE, 1, 1999
PAST MATCH: JUANITA LAKE, 1, 2006 -> JUANITA LAKE, 1, 2005
PAST MATCH: JUANITA LAKE, 2, 2001 -> JUANITA LAKE, 2, 2000
PAST MATCH: JUANITA LAKE, 2, 2002 -> JUANITA LAKE, 2, 2000
PAST MATCH: JUANITA LAKE, 2, 2006 -> JUANITA LAKE, 2, 2005
PAST MATCH: JUANITA LAKE, 3, 1993 -> JUANITA LAKE, 3, 1992
PAST MATCH: JUANITA LAKE, 3, 2001 -> JUANITA LAKE, 3, 2000
PAST MATCH: JUANITA LAKE, 3, 2002 -> JUANITA LAKE, 3, 2000
PAST MATCH: JUANITA LAKE, 3, 2006 -> JUANITA LAKE, 3, 2005
PAST MATCH: JUANITA LAKE, 4, 2001 -> JUANITA LAKE, 4, 2000
PAST MATCH: JUANITA LAKE, 4, 2002 -> JUANITA LAKE, 4, 2000
PAST MATCH: JUANITA LAKE, 4, 2006 -> JUANITA LAKE, 4, 2005
PAST MATCH: JUANITA LAKE, 5, 1999 -> JUANITA LAKE, 5, 1998
PAST MATCH: JUANITA LAKE, 5, 2001 -> JUANITA LAKE, 5, 2000
PAST MATCH: JUANITA LAKE, 5, 2002 -> JUANITA LAKE, 5, 2000
PAST MATCH: JUANITA LAKE, 5, 2005 -> JUANITA LAKE, 5, 2004
PAST MATCH: JUANITA LAKE, 5, 2006 -> JUANITA LAKE, 5, 2004
PAST MATCH: JUANITA LAKE, 6, 1999 -> JUANITA LAKE, 6, 1998
PAST MATCH: JUANITA LAKE, 6, 2001 -> JUANITA LAKE, 6, 2000
PAST MATCH: JUANITA LAKE, 6, 2006 -> JUANITA LAKE, 6, 2005
PAST MATCH: JUANITA LAKE, 7, 2006 -> JUANITA LAKE, 7, 2005
PAST MATCH: JUANITA LAKE, 8, 1999 -> JUANITA LAKE, 8, 1998
PAST MATCH: JUANITA LAKE, 9, 1999 -> JUANITA LAKE, 9, 1998
FUTURE MATCH: JUANITA LAKE, 11, 1992 -> JUANITA LAKE, 11, 1993
PAST MATCH: JUANITA LAKE, 11, 1994 -> JUANITA LAKE, 11, 1993
PAST MATCH: JUANITA LAKE, 12, 1993 -> JUANITA LAKE, 12, 1992
PAST MATCH: JUANITA LAKE, 12, 1994 -> JUANITA LAKE, 12, 1992
PAST MATCH: JUANITA LAKE, 12, 1999 -> JUANITA LAKE, 12, 1998
PAST MATCH: JUANITA LAKE, 12, 2000 -> JUANITA LAKE, 12, 1998
PAST MATCH: JUANITA LAKE, 12, 2001 -> JUANITA LAKE, 12, 1998
PAST MATCH: JUANITA LAKE, 12, 2005 -> JUANITA LAKE, 12, 2004
FUTURE MATCH: EEL RIVER, 1, 1992 -> EEL RIVER, 1, 1993
PAST MATCH: EEL RIVER, 1, 1995 -> EEL RIVER, 1, 1994
PAST MATCH: EEL RIVER, 1, 1997 -> EEL RIVER, 1, 1996
FUTURE MATCH: EEL RIVER, 2, 1992 -> EEL RIVER, 2, 1993
PAST MATCH: EEL RIVER, 2, 1999 -> EEL RIVER, 2, 1998

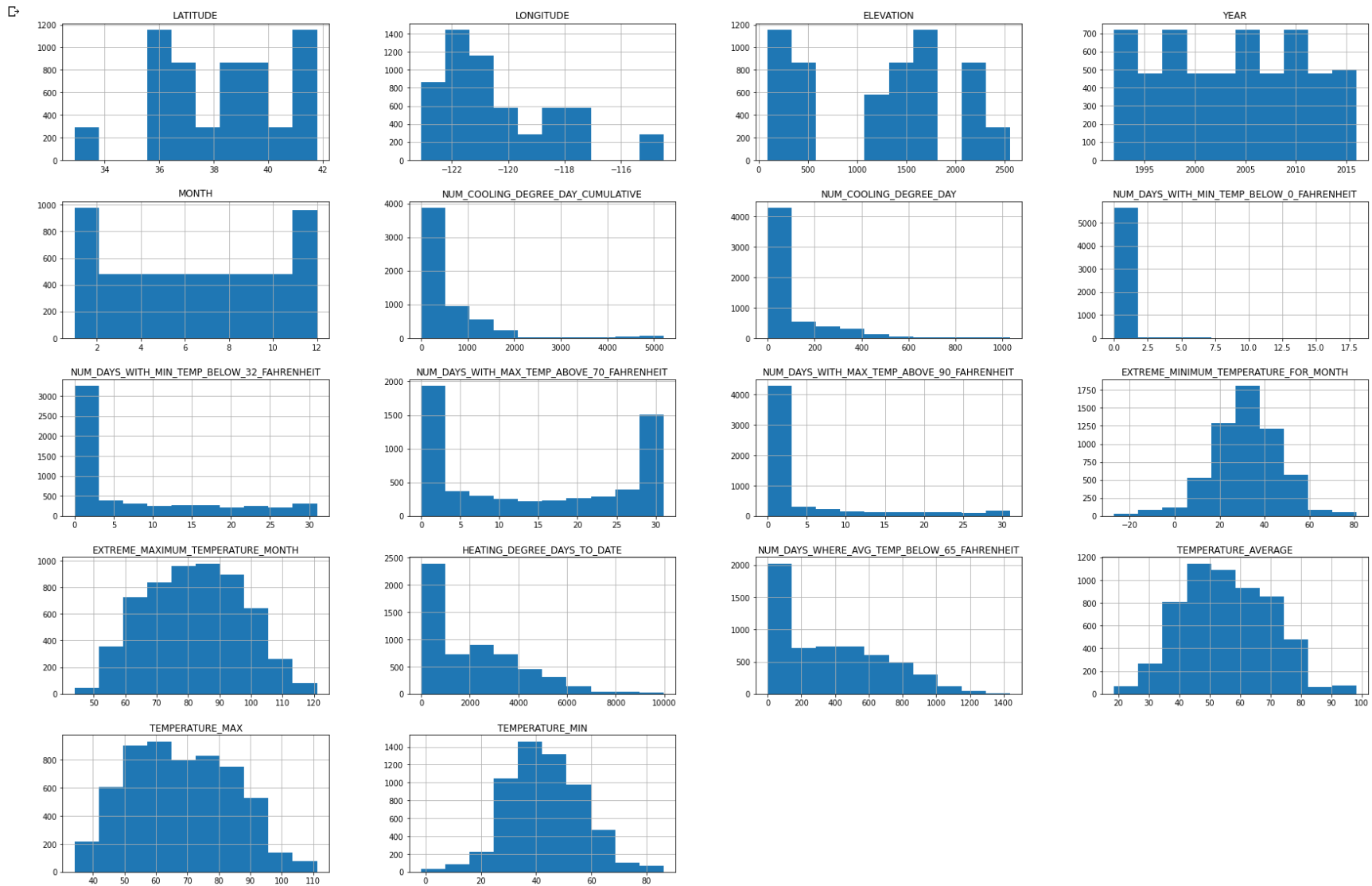
```

#### ▼ Plotting each of the numerical features

```

1 def plot_numerical_features(df):
2     df.select_dtypes(include = "number").hist(figsize=(30, 20))
3
4 plot_numerical_features(noaa_weather)

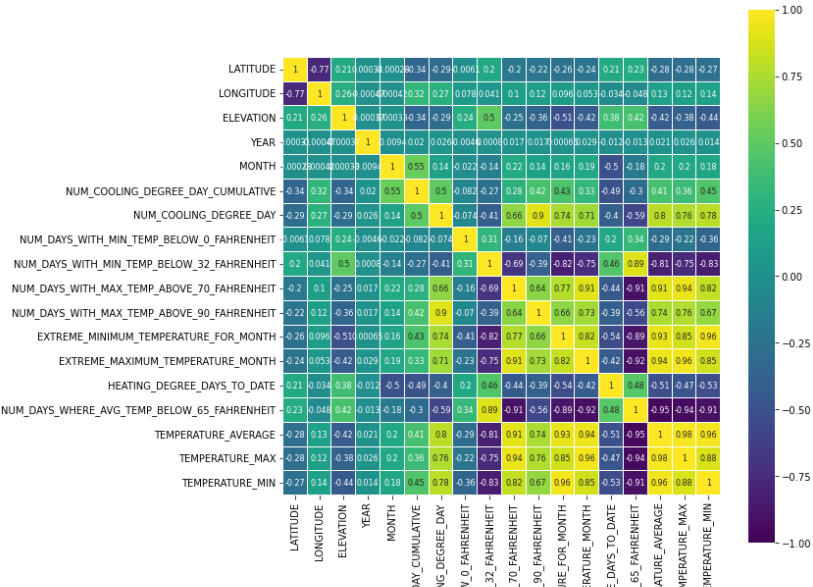
```



```

1 # Correlation Analysis
2 noaa_corr = noaa_weather.select_dtypes(include = "number").corr()
3
4 plt.figure(figsize=(10, 10))
5 sns.heatmap(noaa_corr,
6 cmap='viridis', vmax=1.0, vmin=-1.0, linewidths=0.1,
7 annot=True, annot_kws={"size": 8}, square=True);

```



▼ Correct Skew of Data

```
1 from scipy.stats import skew
2
3 numeric_feats = noaa_weather.dtypes[noaa_weather.dtypes != "object"]
4 numeric_feats = numeric_feats.drop(["LATITUDE", "LONGITUDE", "YEAR", "MONTH"])
5 numeric_feats = numeric_feats.index
6
7 # Check skew in numerical features:
8 skewed_feats = noaa_weather[numeric_feats].apply(lambda x: skew(x)).sort_values(ascending=False)
9 print(skewed_feats)
```

```
NUM_DAYS_WITH_MIN_TEMP_BELOW_0_FAHRENHEIT    9.213557
NUM_COOLING_DEGREE_DAY_CUMULATIVE             3.124595
NUM_COOLING_DEGREE_DAY                       2.682703
NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT    2.044313
NUM_DAYS_WITH_MIN_TEMP_BELOW_32_FAHRENHEIT    1.081981
HEATING_DEGREE_DAYS_TO_DATE                   0.934166
NUM_DAYS_WHERE_AVG_TEMP_BELOW_65_FAHRENHEIT   0.568674
TEMPERATURE_AVERAGE                         0.180799
TEMPERATURE_MAX                             0.168518
NUM_DAYS_WITH_MAX_TEMP_ABOVE_70_FAHRENHEIT    0.122478
TEMPERATURE_MIN                             0.107454
EXTREME_MAXIMUM_TEMPERATURE_MONTH              0.019898
ELEVATION                                    -0.081768
EXTREME_MINIMUM_TEMPERATURE_FOR_MONTH          -0.214279
dtype: float64
```

```
1 from scipy.special import boxcox1p
2
3 # Next, we will fixed highly skewed features.
4 skewed_feats = skewed_feats[abs(skewed_feats) > 0.75]
5 skewed_feats_index = skewed_feats.index
6 lam = 0.15
7 for feat in skewed_feats_index:
8     noaa_weather[feat] = boxcox1p(noaa_weather[feat], lam)
9 print("Corrected skew in the numerical features!")
10
11 plot_numerical_features(noaa_weather)
```



Corrected skew in the numerical features!



## ▼ Normalization

```

1 from sklearn.preprocessing import StandardScaler, MinMaxScaler
2 print(noaa_weather.info())
3 labels_std_normalize = ["ELEVATION", "NUM_COOLING_DEGREE_DAY_CUMULATIVE", "NUM_COOLING_DEGREE_DAY", "NUM_DAYS_WITH_MIN_TEMP_BELOW_0_FAHRENHEIT", "NUM_DAYS_WITH_MIN_TEMP_BELOW_32_FAHRENHEIT", "NUM_DAYS_WITH_MAX_TEMP_ABOVE_70_FAHRENHEIT"]
4 for label in labels_std_normalize:
5     mean = noaa_weather[label].mean()
6     std = noaa_weather[label].std()
7     noaa_weather[label] = (noaa_weather[label] - mean) / std
8 display(noaa_weather)

```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 5779 entries, 0 to 8386
Data columns (total 20 columns):
#   Column                                     Non-Null Count  Dtype
---  -
0    STATION                                     5779 non-null   object
1    LATITUDE                                   5779 non-null   float64
2    LONGITUDE                                  5779 non-null   float64
3    ELEVATION                                  5779 non-null   float64
4    STATION_NAME                               5779 non-null   object
5    YEAR                                       5779 non-null   int64
6    MONTH                                      5779 non-null   int64
7    NUM_COOLING_DEGREE_DAY_CUMULATIVE         5779 non-null   float64
8    NUM_COOLING_DEGREE_DAY                     5779 non-null   float64
9    NUM_DAYS_WITH_MIN_TEMP_BELOW_0_FAHRENHEIT  5779 non-null   float64
10   NUM_DAYS_WITH_MIN_TEMP_BELOW_32_FAHRENHEIT 5779 non-null   float64
11   NUM_DAYS_WITH_MAX_TEMP_ABOVE_70_FAHRENHEIT 5779 non-null   float64
12   NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT 5779 non-null   float64
13   EXTREME_MINIMUM_TEMPERATURE_FOR_MONTH      5779 non-null   float64
14   EXTREME_MAXIMUM_TEMPERATURE_MONTH          5779 non-null   float64
15   HEATING_DEGREE_DAYS_TO_DATE                5779 non-null   float64
16   NUM_DAYS_WHERE_AVG_TEMP_BELOW_65_FAHRENHEIT 5779 non-null   float64
17   TEMPERATURE_AVERAGE                       5779 non-null   float64
18   TEMPERATURE_MAX                            5779 non-null   float64
19   TEMPERATURE_MIN                            5779 non-null   float64
dtypes: float64(16), int64(2), object(2)
memory usage: 948.1+ KB
None
```

	STATION	LATITUDE	LONGITUDE	ELEVATION	STATION_NAME	YEAR	MONTH	NUM_COOLING_DEGREE_DAY_CUMULATIVE	NUM_COOLING_DEGREE_DAY	NUM_DAYS_WITH_MIN_TEMP_BE
0	USR0000CCOH	39.8717	-121.7689	-0.913217	COHASSET	1992	1	-1.239760	-0.886171	
1	USR0000CCOH	39.8717	-121.7689	-0.913217	COHASSET	1992	2	-0.850883	-0.338919	
2	USR0000CCOH	39.8717	-121.7689	-0.913217	COHASSET	1992	3	-0.850883	-0.886171	
3	USR0000CCOH	39.8717	-121.7689	-0.913217	COHASSET	1992	4	-0.316022	0.351139	
4	USR0000CCOH	39.8717	-121.7689	-0.913217	COHASSET	1992	5	0.477971	1.449139	
...	...	...	...	...	...	...	...	...	...	...
8394	USR0000CHNM	36.5625	-117.4736	1.197506	HUNTER_MOUNTAIN	2011	11	0.739154	-0.886171	
8199	USR0000CHNM	36.5625	-117.4736	1.197506	HUNTER_MOUNTAIN	1994	12	0.728818	-0.886171	

Write the data to a CSV for use in later processing

```
8386 USR0000CHNM 36.5625 -117.4736 1.197506 HUNTER_MOUNTAIN 2012 12 0.772949 -0.886171

1 get_matching_station = lambda weather_table, station, month, year: (weather_table[((weather_table["STATION_NAME"] == station) & (weather_table["MONTH"] == month) & (weather_table["YEAR"] == year))])
2 get_matching_station(noaa_weather, "HELL_HOLE", 4, 2004)
3 # HUNTER_MOUNTAIN, 12, 2012
4 get_matching_station(noaa_weather, "HUNTER_MOUNTAIN", 12, 2012)
```

	STATION	LATITUDE	LONGITUDE	ELEVATION	STATION_NAME	YEAR	MONTH	NUM_COOLING_DEGREE_DAY_CUMULATIVE	NUM_COOLING_DEGREE_DAY	NUM_DAYS_WITH_MIN_TEMP_BE
8386	USR0000CHNM	36.5625	-117.4736	1.197506	HUNTER_MOUNTAIN	2012	12	0.772949	-0.886171	

```
1 noaa_weather.to_csv("/content/drive/Shared drives/Data Science 303 Group Project/csv/noaa_weather/noaa_CA_1992_2016_weather_2781174_CLEANED.csv")
```

