

```
!pip install pandas
```

```

Requirement already satisfied: pandas in /usr/local/lib/python3.11/dist-packages
Requirement already satisfied: numpy>=1.23.2 in /usr/local/lib/python3.11/dist-packages
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.11/dist-packages
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packages

```

```

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

```

```

Drive already mounted at /content/drive; to attempt to forcibly remount, call

```

```

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import re
from tqdm import tqdm
import pandas as pd
import numpy as np
import os
import time
import os, time, warnings
from tqdm import tqdm
from multiprocessing import Pool, cpu_count
warnings.filterwarnings("ignore", category=RuntimeWarning)

```

```

pd.set_option('display.max_rows', None)
pd.set_option('display.max_columns', None)
pd.set_option('display.float_format', '{:.2f}'.format)

```

```
df1=pd.read_csv('/content/drive/MyDrive/FE/FE462.csv')
```

```

<ipython-input-21-efb793941f81>:1: DtypeWarning: Columns (14,15,16,17,19,20)
df1=pd.read_csv('/content/drive/MyDrive/FE/FE462.csv')

```

```
df=pd.read_csv('/content/drive/MyDrive/FE/@20.csv')
```

```
total_cells = df.size
missing_cells = df.isnull().sum().sum()
missing_ratio = (missing_cells / total_cells) * 100


print(f" Total cells      : {total_cells:,}")
print(f" Missing cells      : {missing_cells:,}")
print(f" Missing percentage: {missing_ratio:.2f}%")
```

 [Show hidden output](#)

```
missing_report = df.isnull().sum().to_frame(name='Missing Count')
missing_report['Total Rows'] = len(df)
missing_report['Missing %'] = (missing_report['Missing Count'] / missing_report

missing_report = missing_report[missing_report['Missing Count'] > 0]
missing_report = missing_report.sort_values(by='Missing %', ascending=False)

pd.set_option('display.float_format', lambda x: '%.2f' % x)
missing_report
```



	Missing Count	Total Rows	Missing %
<b>Urobilinogen</b>	900000	900000	100.00
<b>Amorphous Elements</b>	900000	900000	100.00
<b>Ketones</b>	900000	900000	100.00
<b>Blood and Haemoglobin</b>	900000	900000	100.00
<b>W.B.Cs / HPF</b>	900000	900000	100.00
<b>Concentration</b>	899996	900000	100.00
<b>Consistancy</b>	899995	900000	100.00
<b>R.B.Cs / HPF</b>	899995	900000	100.00
<b>24 Hour Urine Volume (263)</b>	899993	900000	100.00
<b>Albumin in Urine (263)</b>	899992	900000	100.00
<b>Cystatin C</b>	899992	900000	100.00
<b>W. B. Cs / HPF</b>	899947	900000	99.99
<b>Platelet Count</b>	899908	900000	99.99
<b>Leucocyte esterase</b>	899794	900000	99.98

<b>Epithelial Cells / HPF</b>	899794	900000	99.98
<b>Nitrite</b>	899794	900000	99.98
<b>Colour(Urine Physical Examination)</b>	899794	900000	99.98
<b>Aspect(Urine Physical Examination) Ordinal Encoding</b>	899794	900000	99.98
<b>Protein(Urine Physical Examination)</b>	899794	900000	99.98
<b>Specific Gravity</b>	899794	900000	99.98
<b>Bilirubin_Numeric</b>	899794	900000	99.98
<b>Amorphous_Numeric</b>	899794	900000	99.98
<b>pH(Urine Physical Examination)</b>	899794	900000	99.98
<b>Glucose(Urine Physical Examination)</b>	899794	900000	99.98
<b>Red cell count</b>	899790	900000	99.98
<b>RDW</b>	899787	900000	99.98
<b>MCH</b>	899784	900000	99.98
<b>MCV</b>	899784	900000	99.98
<b>Total Leucocytic Count</b>	899784	900000	99.98
<b>Lymphocytes absolute count</b>	899784	900000	99.98
<b>Basophils absolute count</b>	899784	900000	99.98
<b>Eosinophils absolute count</b>	899784	900000	99.98
<b>MCHC</b>	899784	900000	99.98
<b>Monocytes absolute count</b>	899784	900000	99.98
<b>Neutrophils absolute count</b>	899784	900000	99.98
<b>Hematocrit</b>	899781	900000	99.98
<b>Hemoglobin</b>	899781	900000	99.98
<b>T. Cholesterol/HDL_Numeric</b>	899758	900000	99.97
<b>T. Cholesterol/HDL</b>	899758	900000	99.97
<b>LDL / HDL</b>	899757	900000	99.97
<b>Rheumatoid Factor (quantitative)</b>	899498	900000	99.94
<b>Lead in blood</b>	899290	900000	99.92
<b>Transferrin</b>	898281	900000	99.81

<b>Titre on Hep2 cells</b>	897668	900000	99.74
<b>Microalbuminuria (24 h urine)</b>	890469	900000	98.94
<b>C-Reactive Protein (CRP) quantitative</b>	888923	900000	98.77
<b>Magnesium (Mg) in Serum</b>	882887	900000	98.10
<b>Prostatic Specific Antigen (PSA) Total</b>	780763	900000	86.75
<b>Testosterone (Total)</b>	776737	900000	86.30
<b>Erythrocyte Sedimentation Rate(ESR)</b>	755744	900000	83.97
<b>Ferritin In Serum</b>	749229	900000	83.25
<b>Iron (Fe) in Serum</b>	746343	900000	82.93
<b>Urea in Serum</b>	687717	900000	76.41
<b>BUN/Creatinine Ratio</b>	624730	900000	69.41
<b>Globulin in Serum</b>	624329	900000	69.37
<b>Total Protein in Serum</b>	618897	900000	68.77
<b>Alkaline Phosphatase</b>	611452	900000	67.94
<b>Chloride in Serum</b>	607973	900000	67.55
<b>Diastolic Pressure</b>	607167	900000	67.46
<b>Systolic Pressure</b>	607167	900000	67.46
<b>Estimated Glomerular Filtration Rate(eGFR)</b>	601105	900000	66.79
<b>Sodium (Na) in Serum</b>	597269	900000	66.36
<b>HDL Cholesterol</b>	594795	900000	66.09
<b>Potassium (K) in Serum</b>	593976	900000	66.00
<b>CRP H.S</b>	586365	900000	65.15
<b>Bilirubin (Direct)</b>	581936	900000	64.66
<b>Bilirubin (Total)</b>	581615	900000	64.62
<b>Triglycerides (TG) in Serum</b>	556776	900000	61.86
<b>Cholesterol</b>	551734	900000	61.30
<b>Mean of blood glucose</b>	549685	900000	61.08
<b>Hb A1c %</b>	549093	900000	61.01
<b>Blood Urea Nitrogen (BUN)</b>	538755	900000	59.86
<b>Fasting T4</b>	510510	900000	56.70

<b>Free 14</b>	510513	900000	56.12
<b>Albumin in Serum</b>	468039	900000	52.00
<b>Glucose in Plasma (Fasting)</b>	376147	900000	41.79
<b>Calcium in Serum (Total)</b>	361075	900000	40.12
<b>Thyroid Stimulating Hormone (TSH)</b>	355299	900000	39.48
<b>Uric Acid in Serum</b>	351487	900000	39.05
<b>Vitamin D (25 OH-Vit D -Total)</b>	331532	900000	36.84
<b>Aspartate Aminotransferase (AST)</b>	325714	900000	36.19
<b>Alanine Aminotransferase (ALT)</b>	320871	900000	35.65
<b>Creatinine in Serum</b>	306509	900000	34.06

df.dtypes



0

<b>Unnamed: 0</b>	int64
<b>RESEARCH_ID</b>	object
<b>SAMPLE_ID</b>	object
<b>COLLECTYEAR</b>	int64
<b>REGN_DATE</b>	object
<b>AGE_YEARS</b>	float64
<b>AGE_DAYS</b>	int64
<b>AGE_MONTHS</b>	int64
<b>HEIGHT</b>	int64
<b>WEIGHT</b>	float64
<b>BMI</b>	float64
<b>Thyroid Stimulating Hormone (TSH)</b>	float64
<b>Uric Acid in Serum</b>	float64
<b>Alanine Aminotransferase (ALT)</b>	float64
<b>Ferritin In Serum</b>	float64
<b>Blood Urea Nitrogen (BUN)</b>	float64
<b>Lymphocytes absolute count</b>	float64

<b>R. B. Cs / HPFs</b>	float64
<b>Aspect(Urine Physical Examination) Ordinal Encoding</b>	float64
<b>Eosinophils absolute count</b>	float64
<b>Vitamin D (25 OH-Vit D -Total)</b>	float64
<b>C-Reactive Protein (CRP) quantitative</b>	float64
<b>Transferrin</b>	float64
<b>Red cell count</b>	float64
<b>Basophils absolute count</b>	float64
<b>Crystals(Urine Microscopic Examination :)</b>	float64
<b>Protein(Urine Physical Examination)</b>	float64
<b>Colour(Urine Physical Examination)</b>	float64
<b>Nitrite</b>	float64
<b>LDL Cholesterol</b>	int64
<b>LDL / HDL</b>	float64
<b>24 Hour Urine Volume (263)</b>	float64
<b>Hemoglobin</b>	float64
<b>Total Leucocytic Count</b>	float64
<b>Hematocrit</b>	float64
<b>MCV</b>	float64
<b>Glucose(Urine Physical Examination)</b>	float64
<b>Urea in Serum</b>	float64
<b>Prostatic Specific Antigen (PSA) Total</b>	float64
<b>Testosterone (Total)</b>	float64
<b>Alkaline Phosphatase</b>	float64
<b>Total Protein in Serum</b>	float64
<b>Estimated Glomerular Filtration Rate(eGFR)</b>	float64
<b>Anti CCP Abs</b>	int64
<b>BUN/Creatinine Ratio</b>	float64
<b>Ketones</b>	float64
<b>MCV</b>	float64

	float64
<b>pH(Urine Physical Examination)</b>	float64
<b>Amorphous Elements</b>	float64
<b>Blood and Haemoglobin</b>	float64
<b>Epithelial Cells / HPF</b>	float64
<b>Casts(Urine Microscopic Examination :)</b>	int64
<b>Chloride in Serum</b>	float64
<b>Cholesterol</b>	float64
<b>T. Cholesterol/HDL</b>	float64
<b>Urobilinogen</b>	float64
<b>R.B.Cs / HPF</b>	float64
<b>Erythrocyte Sedimentation Rate(ESR)</b>	float64
<b>Glucose in Plasma (Fasting)</b>	float64
<b>Hb A1c %</b>	float64
<b>Mean of blood glucose</b>	float64
<b>Microalbuminuria (24 h urine)</b>	float64
<b>Bilirubin (Total)</b>	float64
<b>Florescence Pattern</b>	int64
<b>Lead in blood</b>	float64
<b>Monocytes absolute count</b>	float64
<b>Consistancy</b>	float64
<b>Neutrophils absolute count</b>	float64
<b>Specific Gravity</b>	float64
<b>W. B. Cs / HPF</b>	float64
<b>Aspartate Aminotransferase (AST)</b>	float64
<b>Calcium in Serum (Total)</b>	float64
<b>Free T4</b>	float64
<b>Potassium (K) in Serum</b>	float64
<b>Albumin in Serum</b>	float64
<b>Iron (Fe) in Serum</b>	float64

<b>CRP H.S</b>	float64
<b>Triglycerides (TG) in Serum</b>	float64
<b>Rheumatoid Factor (quantitative)</b>	float64
<b>Platelet Count</b>	float64
<b>Albumin in Urine (263)</b>	float64
<b>MCH</b>	float64
<b>RDW</b>	float64
<b>W.B.Cs / HPF</b>	float64
<b>Leucocyte esterase</b>	float64
<b>Concentration</b>	float64
<b>Creatinine in Serum</b>	float64
<b>Sodium (Na) in Serum</b>	float64
<b>Bilirubin (Direct)</b>	float64
<b>Magnesium (Mg) in Serum</b>	float64
<b>Titre on Hep2 cells</b>	float64
<b>HDL Cholesterol</b>	float64
<b>Globulin in Serum</b>	float64
<b>Cystatin C</b>	float64
<b>RESEARCH_ID_int</b>	int64
<b>GENDER_BINARY</b>	int64
<b>CITY_NAME_ENCODED</b>	int64
<b>BDL</b>	int64
<b>Florescence Pattern</b>	int64
<b>Systolic Pressure</b>	float64
<b>Diastolic Pressure</b>	float64
<b>Amorphous_Numeric</b>	float64
<b>Bilirubin_Numeric</b>	float64
<b>T. Cholesterol/HDL_Numeric</b>	float64

**dtype:** object



```
object_columns_df = df.select_dtypes(include=['object'])
object_columns_df.dtypes
```

```
0
RESEARCH_ID  object
SAMPLE_ID    object
REGN_DATE    object
```

**dtype:** object

```
df.duplicated().sum()
```

```
np.int64(0)
```

```
df.shape
```

```
(900000, 104)
```

```
df.columns
```

```
Index(['Unnamed: 0', 'RESEARCH_ID', 'SAMPLE_ID', 'COLLECTYEAR',
      'REGN_DATE',
      'AGE_YEARS', 'AGE_DAYS', 'AGE_MONTHS', 'HEIGHT', 'WEIGHT',
      ...,
      'RESEARCH_ID_int', 'GENDER_BINARY', 'CITY_NAME_ENCODED', 'BDL',
      'Florescence Pattern ', 'Systolic Pressure', 'Diastolic Pressure',
      'Amorphous_Numeric', 'Bilirubin_Numeric', 'T.
Cholesterol/HDL_Numeric'],
      dtype='object', length=104)
```

```
columns_zero_is_missing = [
    "LDL Cholesterol",
    "HDL Cholesterol",
    "LDL / HDL",
    "Triglycerides (TG) in Serum",
    "Cholesterol",
    "T. Cholesterol/HDL",
    "T. Cholesterol/HDL_Numeric",
    "Bilirubin (Total)",
    "Bilirubin (Direct)",
    "Iron (Fe) in Serum",
    "Ferritin In Serum",
    "Uric Acid in Serum",
```

```

"Albumin in Serum",
"Globulin in Serum",
"Total Protein in Serum",
"Calcium in Serum (Total)",
"Magnesium (Mg) in Serum",
"Sodium (Na) in Serum",
"Potassium (K) in Serum",
"Creatinine in Serum",
"Urea in Serum",
"Blood Urea Nitrogen (BUN)",
"Estimated Glomerular Filtration Rate(eGFR)",
" BUN/Creatinine Ratio",
"Cystatin C",
"Thyroid Stimulating Hormone (TSH)",
"Free T4",
"Prostatic Specific Antigen (PSA) Total",
"Testosterone (Total)",
"Vitamin D (25 OH-Vit D -Total)",
"C-Reactive Protein (CRP) quantitative",
"CRP H.S",
"Rheumatoid Factor (quantitative)",
"Anti CCP Abs",
"Transferrin",
"Aspartate Aminotransferase (AST)",
"Alanine Aminotransferase (ALT)",
"Alkaline Phosphatase",
"Hemoglobin",
"Total Leucocytic Count",
"Platelet Count",
"MCV",
"MCH",
"MCHC",
"RDW",
"Mean of blood glucose",
"Hb A1c %",
"Glucose in Plasma (Fasting)",
"Microalbuminuria (24 h urine)",
"Albumin in Urine (263)",
"Lead in blood",
"Casts(Urine Microscopic Examination :)",
]

```

```

for col in columns_zero_is_missing:
    if col in df.columns:
        df[col] = df[col].replace(0, np.nan)

```

```
df.isnull().sum().sum()
```

```
df.duplicated().sum().sum()
```

```
df_sorted = df.sort_values(by=["RESEARCH_ID", "SAMPLE_ID", "REGN_DATE"])
df_fast = df_sorted.copy()

columns_to_fill = df_fast.columns[df_fast.isnull().any()].tolist()

total_filled = 0
start_time = time.time()

for col in tqdm(columns_to_fill, desc="Filling NaNs", unit="column"):
    before = df_fast[col].isna().sum()
    df_fast[col] = df_fast.groupby(["RESEARCH_ID", "SAMPLE_ID"])[col].ffill()
    df_fast[col] = df_fast.groupby(["RESEARCH_ID", "SAMPLE_ID"])[col].bfill()
    after = df_fast[col].isna().sum()
    total_filled += (before - after)

end_time = time.time()
print(f"\n✅ Total filled values: {total_filled:,}")
print(f"🕒 Time elapsed: {end_time - start_time:.2f} seconds.")

output_folder = "/content/drive/MyDrive/XXXX/knowledge_project"
os.makedirs(output_folder, exist_ok=True)

output_path = os.path.join(output_folder, "cleaned_dataset.csv")
df_fast.to_csv(output_path, index=False)

print(f"\n File saved successfully at: {output_path}")
```

```
df_fast=pd.read_csv('/content/drive/MyDrive/XXXX/knowledge_project/cleaned_data
```

```
df_fast.duplicated().sum().sum()
```

```
↔ np.int64(0)
```

```
df_fast.isnull().sum().sum()
```

```
np.int64(61344978)
```

```
df_fast["REGN_DATE"] = pd.to_datetime(df_fast["REGN_DATE"], errors='coerce')
```

```
df_fast["REGN_DATE"].dtype
```

```
dtype('<M8[ns]')
```

```
df_fast = df_fast.sort_values(by="RESEARCH_ID")
```

```
df_marg=df_fast
```

```
df_marg.isnull().sum().sum()
```

```
np.int64(61344978)
```

```
def smart_impute_group_numpy_rid(group_tuple):
    _, group = group_tuple
    group = group.copy()
    dates = group["REGN_DATE"].values.astype("datetime64[ns]")
    excluded_cols = ["RESEARCH_ID", "SAMPLE_ID", "REGN_DATE"]
    columns_to_fill = [col for col in group.columns if col not in excluded_cols]

    for col in columns_to_fill:
        values = group[col].values
        mask_nan = np.isnan(values)
        if not np.any(mask_nan): continue
        known_idx = np.where(~mask_nan)[0]
        known_dates = dates[known_idx]
        known_values = values[known_idx]

        for idx in np.where(mask_nan)[0]:
            current_date = dates[idx]
            same_date_mask = known_dates == current_date
            if np.any(same_date_mask):
                values[idx] = known_values[same_date_mask][0]
                continue
            if known_dates.size > 0:
                time_deltas = np.abs(known_dates - current_date)
                values[idx] = known_values[np.argmin(time_deltas)]
                continue
            if known_values.size > 0:
```

```

        mean_val = np.nanmean(known_values)
        if not np.isnan(mean_val):
            values[idx] = mean_val
    group[col] = values
    return group

all_groups = list(df_marg.groupby("RESEARCH_ID"))
batch_size = 50000
total_batches = (len(all_groups) + batch_size - 1) // batch_size
output_folder = "/content/drive/MyDrive/XXXX/rid_batches"
os.makedirs(output_folder, exist_ok=True)

start_time = time.time()

for i in range(total_batches):
    batch_path = f"{output_folder}/batch_rid_{i+1}.csv"
    if os.path.exists(batch_path):
        print(f"✅ Skipping batch {i+1} (already processed).")
        continue

    batch_groups = all_groups[i * batch_size : (i+1) * batch_size]
    results = []

    with Pool(cpu_count() - 1) as pool:
        for result in tqdm(pool.imap(smart_impute_group_numpy_rid, batch_groups),
                           total=len(batch_groups),
                           desc=f"Batch RID {i+1}",
                           unit="group",
                           dynamic_ncols=True):
            results.append(result)

    df_batch = pd.concat(results, ignore_index=True)
    df_batch.to_csv(batch_path, index=False)
    print(f" Saved: {batch_path}")

print("\n Merging all RID batches...")
merged_df = pd.concat(
    [pd.read_csv(f"{output_folder}/batch_rid_{i+1}.csv") for i in range(total_batches)],
    ignore_index=True
)

final_path = "/content/drive/MyDrive/XXXX/filled_final_by_rid_batches.csv"
merged_df.to_csv(final_path, index=False)

```


```
end_time = time.time()
minutes, seconds = divmod(int(end_time - start_time), 60)

print(f" RID Imputation Completed in {total_batches} Batches")
print(f" Final file saved at: {final_path}")
print(f" Remaining missing values: {merged_df.isnull().sum().sum():,}")
print(f"⌚ Time taken: {minutes} min {seconds} sec")
```


 [Show hidden output](#)

```
df4=pd.read_csv('/content/drive/MyDrive/XXXX/filled_final_by_rid_batches.csv')
```

```
df4.isnull().sum().sum()
```

 np.int64(55841104)

```
df4.shape
```


 (900000, 104)

```
df_marg2=pd.read_csv("/content/drive/MyDrive/XXXX/filled_final_by_rid_batches.c
```

```
df_marg2.isnull().sum().sum()
```

 np.int64(55841104)

```
df_marg2.duplicated().sum().sum()
```

 np.int64(0)

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

