

```
[39] ✓ Os
    # Fill missing PED_ROLE values with 'UNKNOWN'
    df_persons['PED_ROLE'] = df_persons['PED_ROLE'].fillna('UNKNOWN')

    print("Missing PED_ROLE values have been imputed with 'UNKNOWN'.")

    Missing PED_ROLE values have been imputed with 'UNKNOWN'.

[40] ✓ Os
    ⏎ # Count missing values
    missing_count = df_persons['PED_ROLE'].isnull().sum()

    # Calculate percentage
    missing_percent = (missing_count / len(df_persons)) * 100

    print(f"Missing values in PED_ROLE: {missing_count} rows")
    print(f"Percentage missing: {missing_percent:.2f}%")

    ... Missing values in PED_ROLE: 0 rows
    Percentage missing: 0.00%
```

**Fill missing PED\_ROLE values with "UNKNOWN" to avoid bias toward the dominant role. This preserves the true distribution of categories. It also keeps imputed records traceable for analysis.**

```
[42] ✓ 1s
    ⏎ # Fill missing PERSON_SEX values with 'U'
    df_persons['PERSON_SEX'] = df_persons['PERSON_SEX'].fillna('U')

    print("Missing PERSON_SEX values have been imputed with 'U'.")

    ... Missing PERSON_SEX values have been imputed with 'U'.

[43] ✓ 0s
    ⏎ # Count missing values
    missing_count = df_persons['PERSON_SEX'].isnull().sum()

    # Calculate percentage
    missing_percent = (missing_count / len(df_persons)) * 100

    print(f"Missing values in PERSON_SEX: {missing_count} rows")
    print(f"Percentage missing: {missing_percent:.2f}%")

    ... Missing values in PERSON_SEX: 0 rows
    Percentage missing: 0.00%
```

**We filled missing PERSON\_SEX values with "U" to match the dataset's existing coding scheme.**

**This avoids bias toward male or female categories during analysis. It also keeps imputed records traceable and consistent for reproducibility.**

••• Այստեղ ենք լախութեզ միւն սex-սեռություն ազդես? - Ե կեր էօւ լուսավոր ենք.

Եւյս („Այստեղ ենք լախութեզ միւն սex-սեռություն ազդես? - Ե կեր էօւ լուսավոր ենք.“)

• բառաչորակ(լախութե՞րլ-սex)  
գէ-ներսոնս[, նԵՐՏՈՒՎԵԵ, ] = գէ-ներսոնս·Բրոներլ(, նԵՐՏՈՒՎԵԵ-SEX, )[. նԵՐՏՈՒՎԵԵ, ]/  
  
Ներուս ՏԵՐԵՍ. ԳՐԱԿԱՐԱ(ազդես-ՎԵՐ). ԱՏՔԼԱՅ(լուֆ)  
ազդես-ՎԵՐ = ՏԵՐԵՍ[ՏԵՐԵՍ] = -Ե. ազդես()  
զեւ լախութե՞րլ-սex(ՏԵՐԵՍ):  
# Լախութե այստեղ լայնը միւն սex-սեռություն ազդես

)  
յամրգ x: x լէ 0 <= x <= 100 Ելց -Ե  
գէ-ներսոնս[, նԵՐՏՈՒՎԵԵ, ] = գէ-ներսոնս[, նԵՐՏՈՒՎԵԵ, ]. ԳԵՐԵԼ(  
# Կերպար լուսավոր ենք միւն -Ե

◀ ▶ [48]

**Imputing missing ages with the median by PERSON\_SEX uses a logical demographic link. It avoids bias from unrealistic values while keeping distributions realistic. This ensures consistency and transparency in the dataset.**

**Unrealistic error outliers (ages <0 or >100) were replaced with -1 as a sentinel marker. This ensures missing values are filled logically while invalid entries remain traceable for analysis.**

```
[51] ✓ 2s ⏎ # Step 1: Find the mode (most frequent value)
mode_value = df_persons['EJECTION'].mode()[0]

# Step 2: Fill missing values with the mode
df_persons['EJECTION'] = df_persons['EJECTION'].fillna(mode_value)

print(f"Missing EJECTION values imputed with mode: {mode_value}")

...
... Missing EJECTION values imputed with mode: Not Ejected
```

  

```
[52] ✓ 0s ⏎ # Count missing values
missing_count = df_persons['EJECTION'].isnull().sum()

# Calculate percentage
missing_percent = (missing_count / len(df_persons)) * 100

print(f"Missing values in EJECTION: {missing_count} rows")
print(f"Percentage missing: {missing_percent:.2f}%")

...
... Missing values in EJECTION: 0 rows
Percentage missing: 0.00%
```

**imputed missing EJECTION values using the mode because it is a categorical variable. The mode represents the most common in this case 97% and therefore most likely category, making it a logical default. This ensures consistency in the dataset while avoiding arbitrary or unrealistic replacements.**

```
... МИЛТИСОГ 9АСУЛГЕ ВАЛЕНЧИСГЕ СЕУЛДИФЕД ВИДИ, ОУИКНОМИ, ''  
    балтнф("МИЛТИСОГ 9АСУЛГЕ ВАЛЕНЧИСГЕ ВАЛДИФЕД ВИДИ, ОУИКНОМИ, '')  
    дт-бенсона[ ,9АСУЛГЕ, ] = дт-бенсона[ ,9АСУЛГЕ, ].титт( "ОУИКНОМИ")  
    # ЭДИ МИЛТИСОГ 9АСУЛГЕ ВАЛЕНЧИСГЕ ВАЛДИФЕД ВИДИ, ОУИКНОМИ, ''
```

We imputed missing POSITION\_IN\_VEHICLE values with "UNKNOWN" because the variable is categorical. Using "UNKNOWN" avoids introducing bias from mode imputation when the true position cannot be determined. It provides a clear, explicit placeholder that keeps the dataset consistent and traceable for later analysis.