

## I. Differences Between the New and Old Datasets

### 1. Missing Columns in the New Dataset:

- *tersiyerhp*
- *Medeni Durum (marital status)*
- *24 saat idrar sodyum (urinary\_sodium\_24h)*
- *Kemik spesifik ALP (bone\_specific\_alp)*
- *2nd and 3rd AMELİYAT PATOLOJİ RAPORU (2nd and 3rd surgery\_pathology\_report)*
- *hdl*

### 2. Renaming of Columns:

- The column formerly named "**nhph(normokals)**" is now renamed to "**NORMOKALSEM?K**".

### 3. New Column Inclusions:

- The new dataset includes "**SEMTOMATIKPHP**", a column that was not part of the previous mapping.
  - Two versions of the ALP measurement exist: "**alp**" and "**ALP**".
- 

## II. Data Preprocessing Steps

### 1. Column Removal:

- **After ALP:** All columns following the ALP column were removed.
- **Sparse Columns:** The following columns were removed due to sparse data:
  - Unnamed: 89, Unnamed: 92, Unnamed: 79, Unnamed: 84, Unnamed: 85, Unnamed: 95, Unnamed: 95.1, Unnamed: 96, Unnamed: 97, Unnamed: 98, Unnamed: 99, Unnamed: 100.

- **Additional Dropped Columns:** The following columns were dropped:
  - 'hospital\_name', 'living\_city', 'mobile\_phone\_number', 'patient\_name\_1st\_visit', 'date\_of\_1st\_visit', 'patient\_name\_2nd\_visit', 'date\_of\_2nd\_visit', 'patient\_name\_3rd\_visit', 'date\_of\_3rd\_visit', 'medication\_name', 'comorbidity', and 'semtomatikphp'.
- The "**semtomatikphp**" target column was removed entirely because it contained no data.

## 2. Handling Dual ALP Columns:

- For the two ALP columns ("**alp**" and "**ALP**"), is following strategy correct:
  - **If one column is empty** and the other contains a value, use the non-empty value.
  - **If both columns have values**, compute the mean of the two values and use that as the final ALP value.

## 3. Categorical Missing Values:

- Should empty cells be filled with mean?

## 4. Conversion of Categorical Variables:

- **aphp Column:** Map 'Asemtomatik' and 'EVET' to 1; missing values are set to 0.
- **nhph Column:** Map 'nhph(normokals)' and 'Normokalsemik' to 1; map 'HAYIR' and missing values to 0.
- **secondary\_hp Column:** Map various affirmative responses (e.g., 'EVET' and entries resembling 'Sekonder hp') to 1; map 'HAYIR' and missing values to 0.

## 5. Gender Encoding:

- 'Men' → 0
- 'Women' → 1
- Missing or invalid values → -1

---

## III. Data Format Corrections

- Empty/Whitespace: " ", " ", " "
- Comma as Decimal: "0,6", "0,69"
- Square Brackets: "[1.0726]", "[103.9000]"
- Multiple Values (semicolon): "8.5; 11.2", "47; 22", "30.6; 115.3"
- Multiple Dots: "8.511.2"
- Thousand Separators: "1,063,125"

## Example Input & Output:

Raw Data ( bmi )	Processed Data ( bmi )
"23.5"	23.5
"22,8"	22.8
"25.0; 26.3"	26.3
"Mg: 10.2"	10.2
"[15.6]"	15.6
".."	0

### 1. Undecided Handling for Specific Diagnostic/Imaging Columns:

- The following columns require further decisions on handling:
  - 'parathyroid\_scintigraphy'
  - 'parathyroid\_spect'
  - 'neck\_ultrasound'
  - '4d\_ct'
  - '4d\_mri'

- '99mtc\_mibi'
- '99mtc\_sestamibi'
- 'first\_surgery\_pathology\_report'