

Pusula Data Science Intern Case

EDA & Preprocessing — Extras

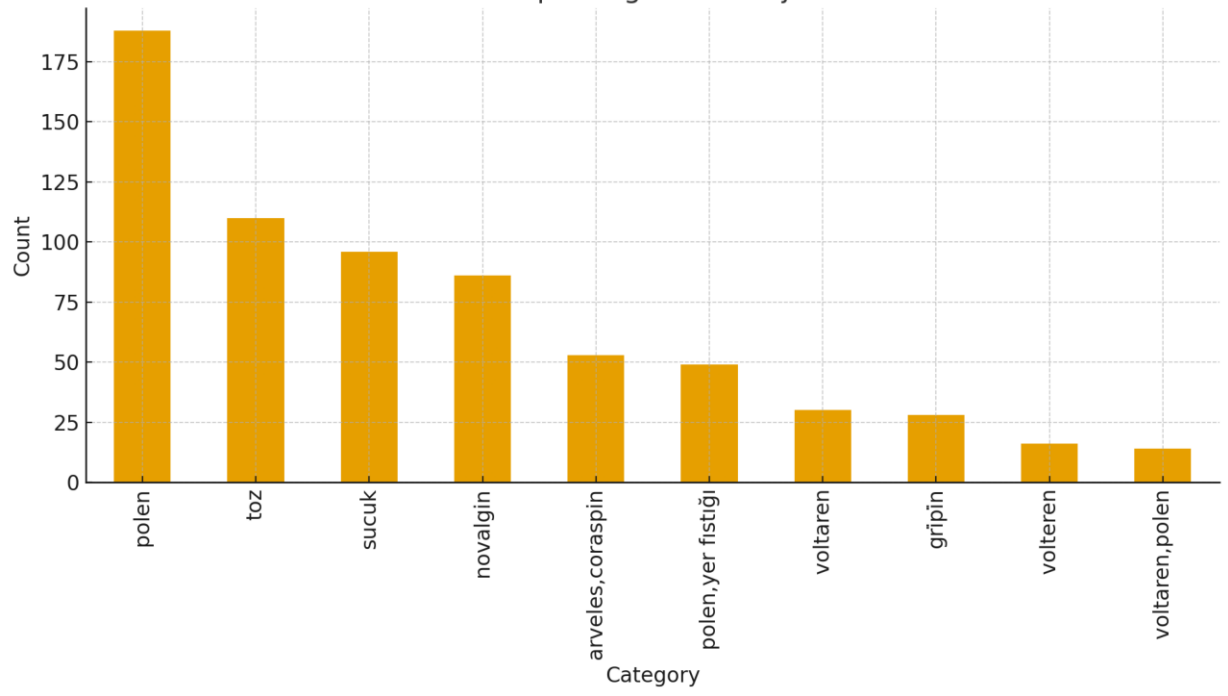
Rows: 1273 | Cols: 13

Key EDA Findings:

- Numeric columns: HastaNo, Yas
- Categorical/text columns: Cinsiyet, KanGrubu, Uyruk, KronikHastalik, Bolum, Alerji, Tanilar, TedaviAdi, TedaviSuresi, UygulamaYerleri, UygulamaSuresi
- Missing present in: Cinsiyet, KanGrubu, KronikHastalik, Bolum, Alerji, Tanilar, UygulamaYerleri
- ID `HastaNo` excluded from features to prevent leakage.
- Top-20 token binarization used for multi-label list columns.

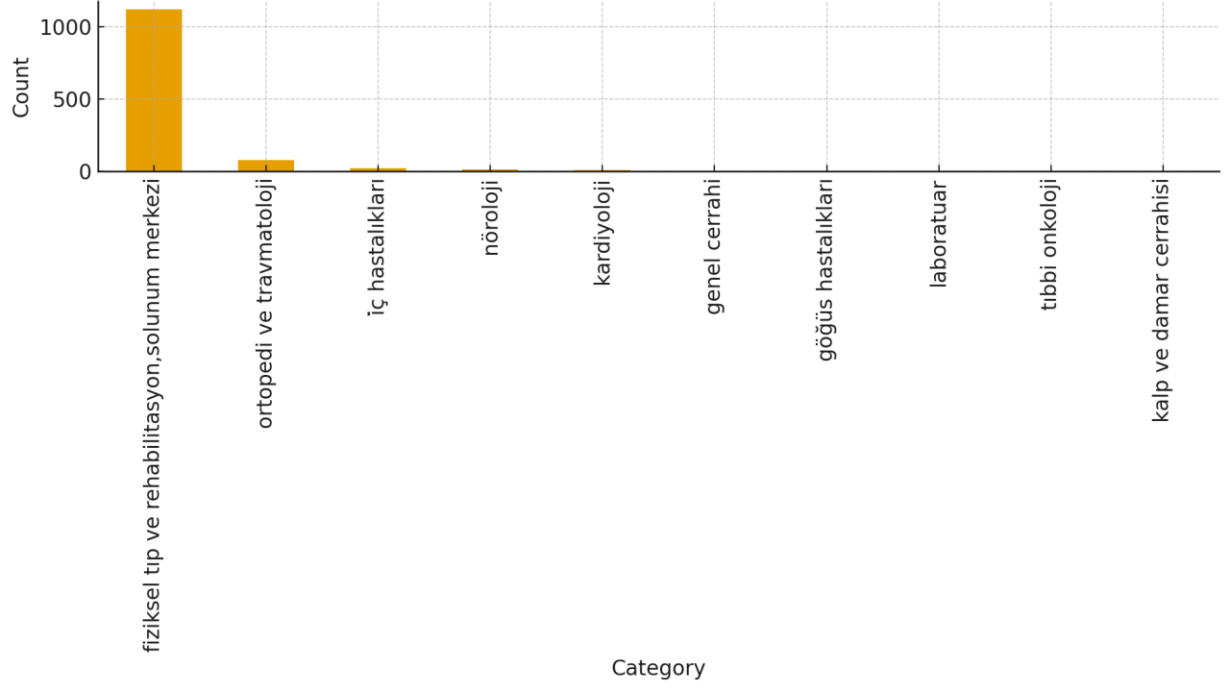
bar_top_Alerji.png

Top categories: Alerji



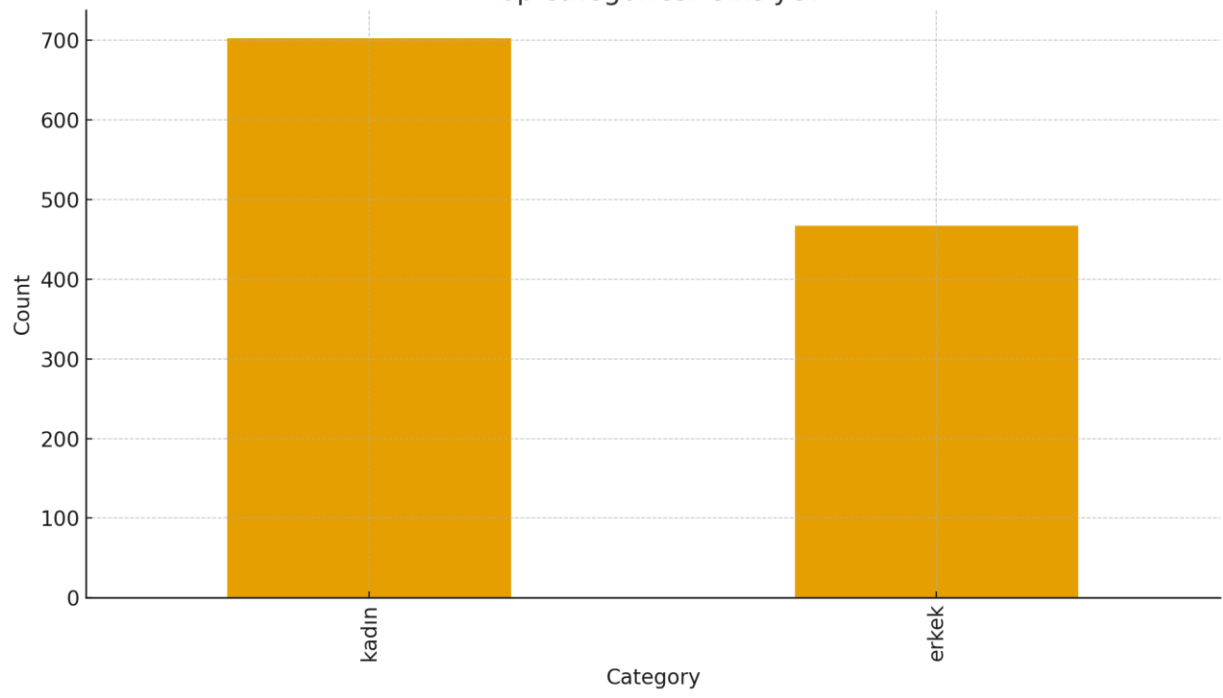
bar_top_Bolum.png

Top categories: Bolum



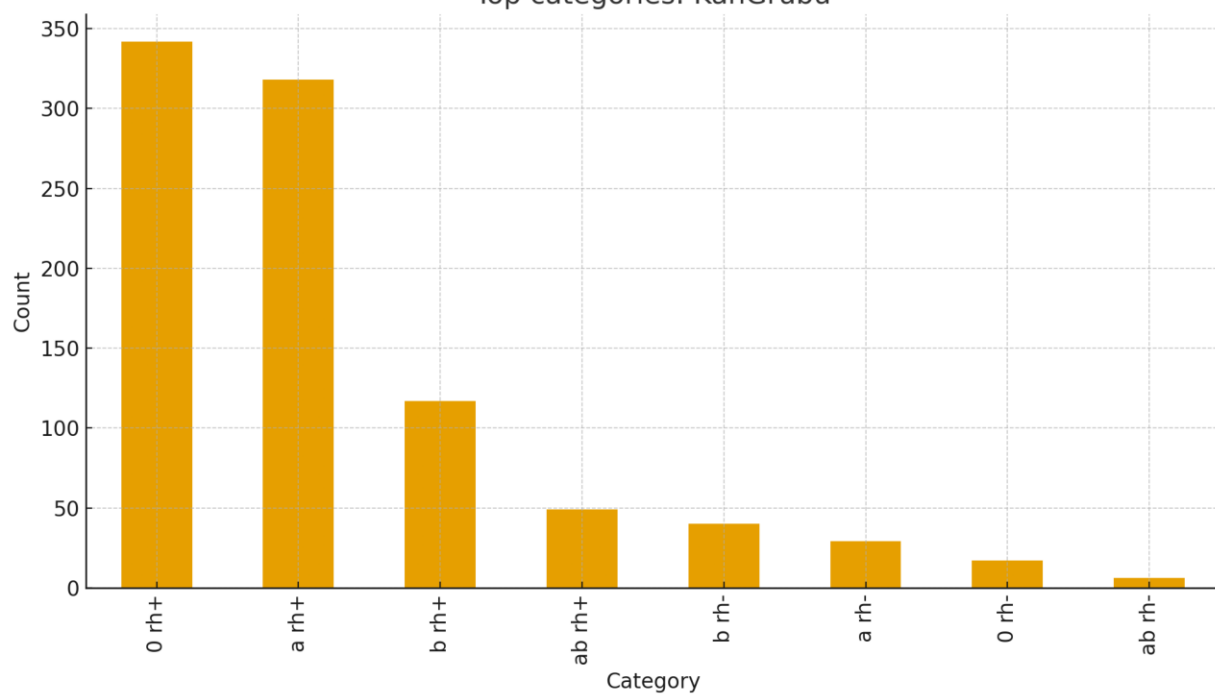
bar_top_Cinsiyet.png

Top categories: Cinsiyet



bar_top_KanGrubu.png

Top categories: KanGrubu



Extras Implemented:

- Pipeline-level, modular preprocessing with ColumnTransformer.
- High-cardinality categorical handling via frequency encoding.
- Multi-label splitter (comma-separated) with Top-K binarization per field.
- Exported cleaned dataset, prepared features, target, and fitted pipeline.
- Matplotlib-only charts (one per figure) for numeric distributions and top categories.
- README & requirements for quick reproducibility.
- Ethics/Privacy: patient ID excluded from features; target-only usage for `TedaviSuresi`.