

PhenoQC Quality Control Report

Source file: synthetic_multi_data.csv

Summary

Imputation Strategy	Mean
Schema Validation Score	0.00%
Missing Data Score	93.52%
Mapping Success Score	99.79%
Overall Quality Score	64.44%

Imputation Settings

Global Strategy	knn
Global Params	{'n_neighbors': 5, 'weights': 'uniform'}
Tuning Enabled	True
Best Params	{'n_neighbors': 7}
Tuning Score	198.6406 (MAE)

Data Quality Scores:

Schema Validation Score: 0.00%
Missing Data Score: 93.52%
Mapping Success Score: 99.79%
Overall Quality Score: 64.44%

Schema Validation Results

Format Validation: False
Duplicate Records: 300 issues found.
Conflicting Records: 300 issues found.
Integrity Issues: 3150 issues found.
Referential Integrity Issues: No issues found.
Anomalies Detected: 44 issues found.

Invalid Mask: 3150 issues found.

Missing Data Summary

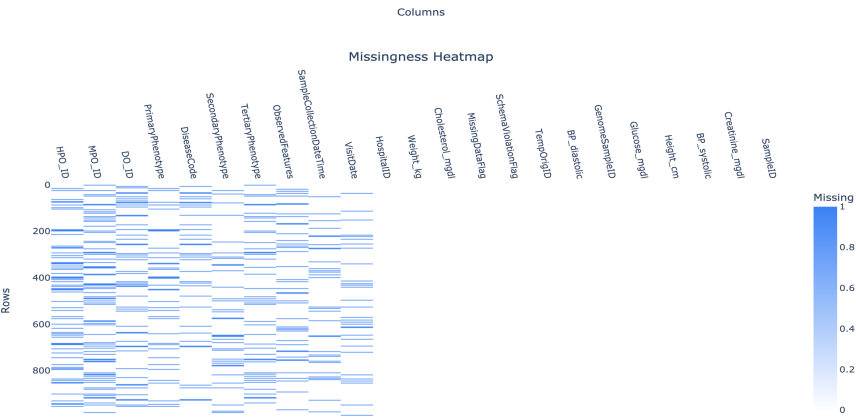
Column	Missing Count
TertiaryPhenotype	397
SecondaryPhenotype	390
PrimaryPhenotype	383
DiseaseCode	364
ObservedFeatures	354
Glucose_mgdl	322
BP_systolic	321
Height_cm	320
BP_diastolic	318
VisitDate	317
Weight_kg	315
Creatinine_mgdl	307
Cholesterol_mgdl	299
SampleCollectionDateTime	288

Records Flagged for Missing Data: 2489

Ontology Mapping Summary

Ontology	Total Terms	Mapped	Success Rate
HPO	1994	1989	99.75%
DO	2066	2062	99.81%
MPO	2030	2026	99.80%

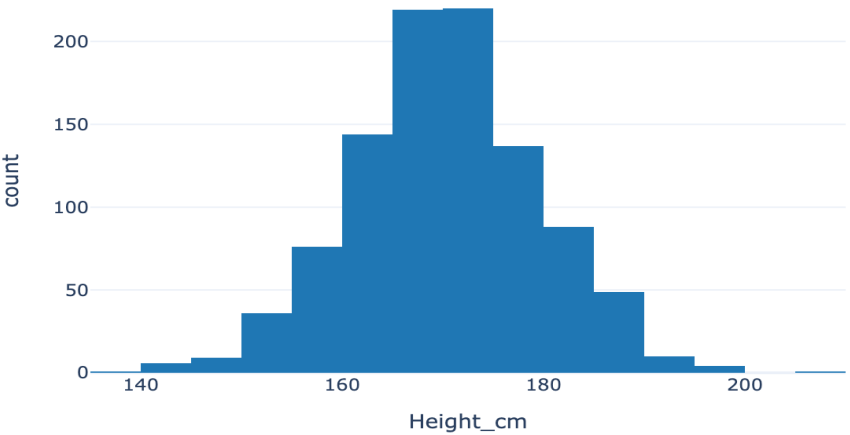
Visualizations



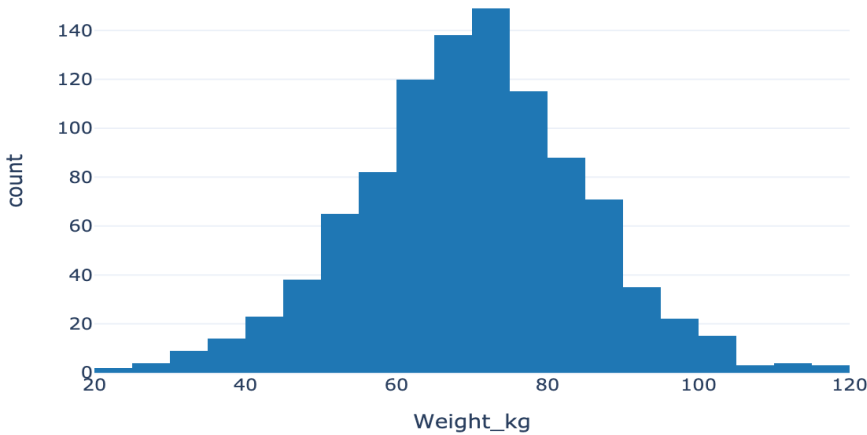
Percentage of Missing Data by Column



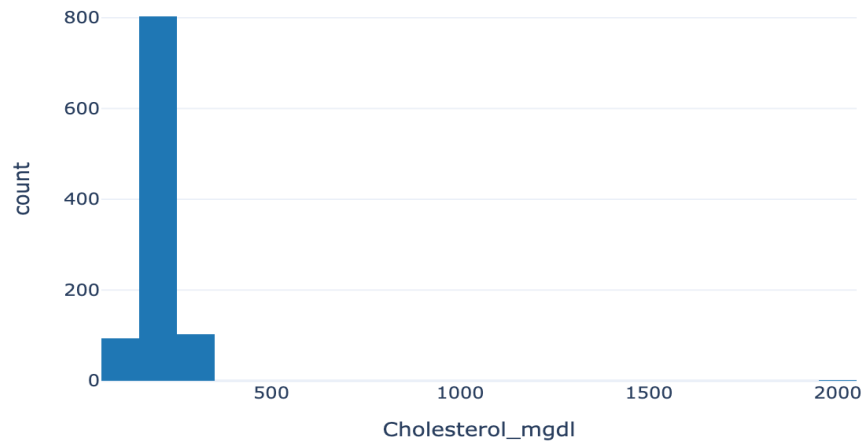
Distribution of Height_cm



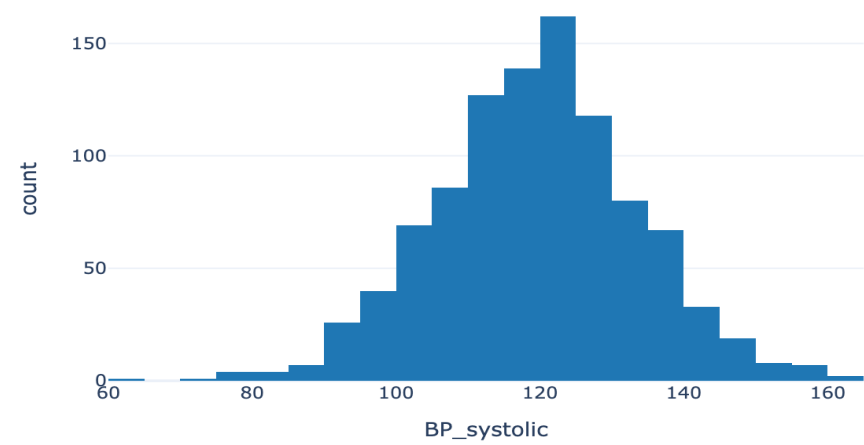
Distribution of Weight_kg



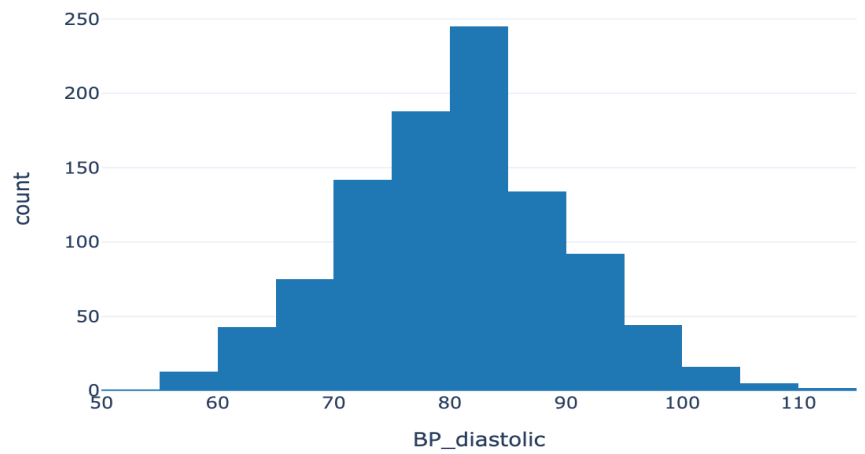
Distribution of Cholesterol_mgdl



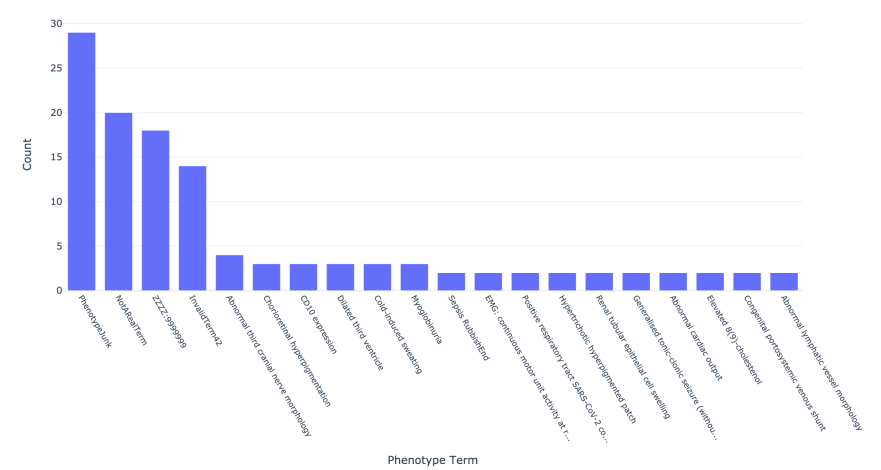
Distribution of BP_systolic



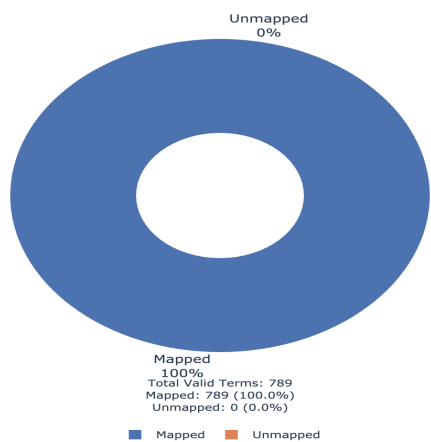
Distribution of BP_diastolic



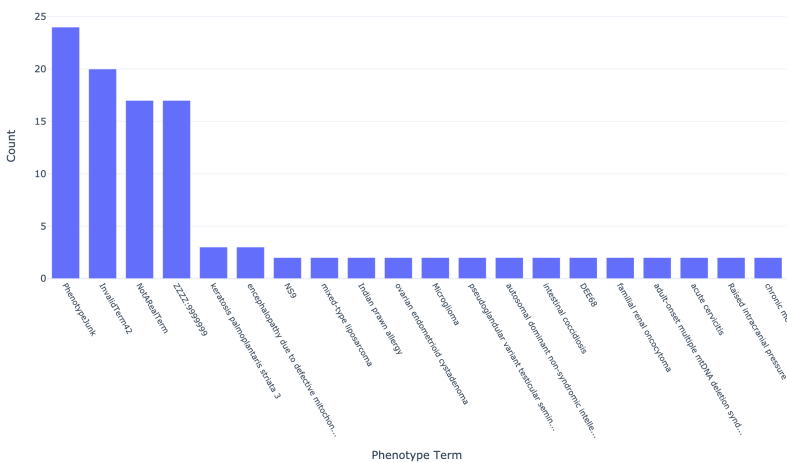
Top 20 Most Common Terms in PrimaryPhenotype



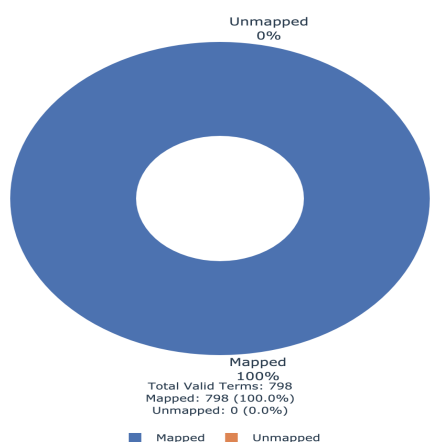
Mapping Results: PrimaryPhenotype → HPO



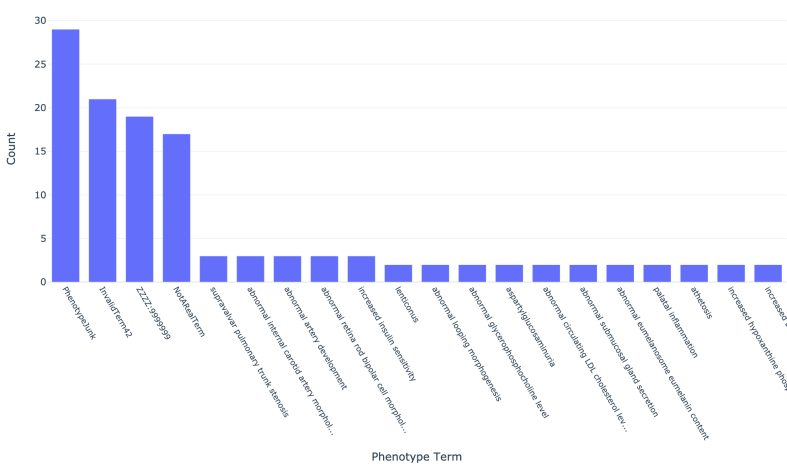
Top 20 Most Common Terms in DiseaseCode



Mapping Results: DiseaseCode → DO



Top 20 Most Common Terms in TertiaryPhenotype



Mapping Results: TertiaryPhenotype → MPO

