

PhenoQC Quality Control Report

Source file: synthetic_multi2_data.csv

Summary

Imputation Strategy	Mean
Schema Validation Score	0.00%
Missing Data Score	92.94%
Mapping Success Score	99.75%
Overall Quality Score	64.23%

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: 92.94%
Mapping Success Score: 99.75%
Overall Quality Score: 64.23%

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

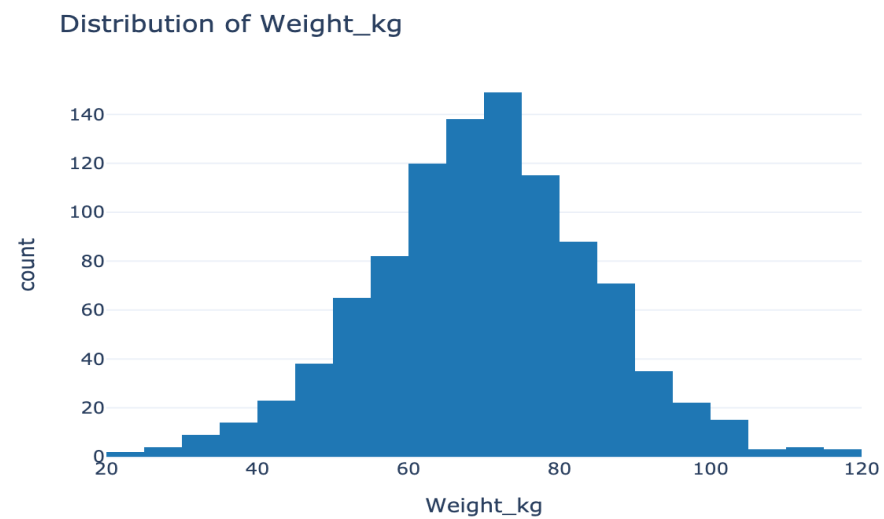
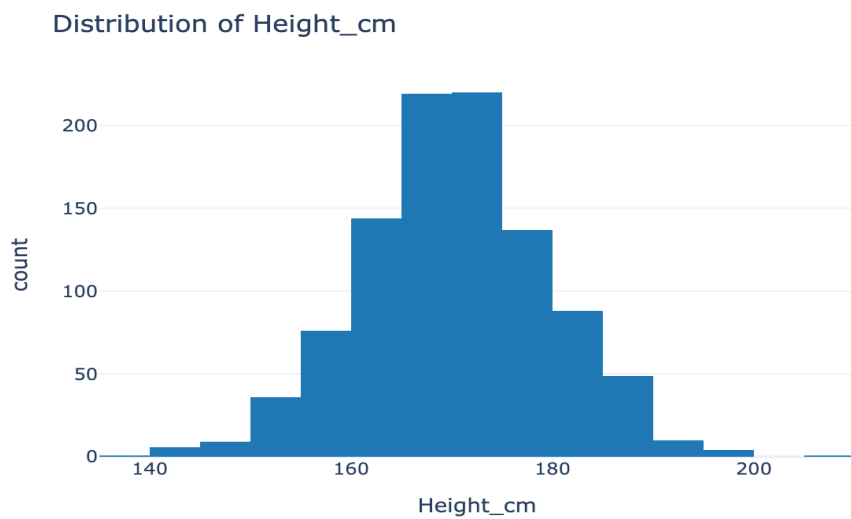
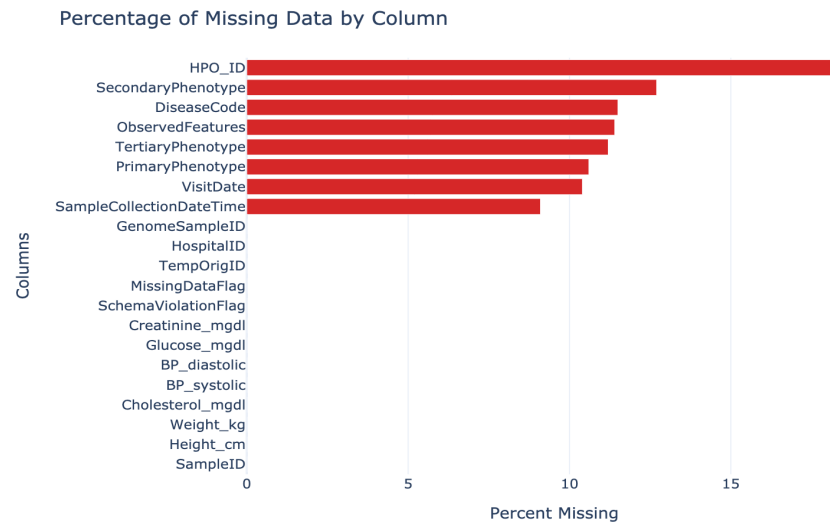
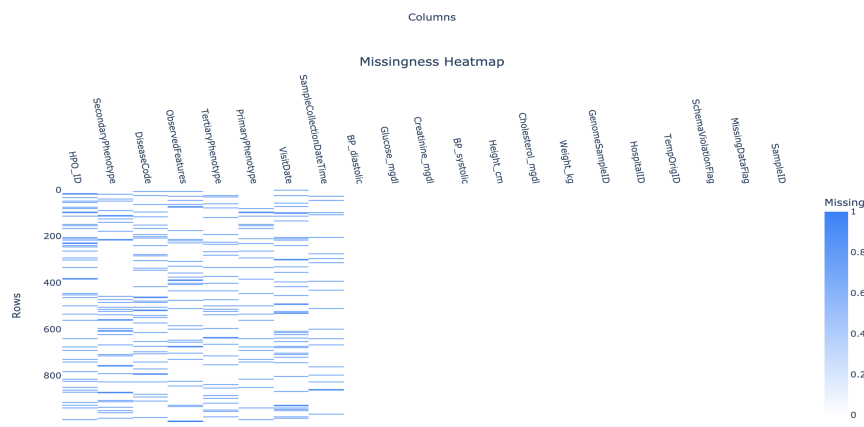
Column	Missing Count
SecondaryPhenotype	399
DiseaseCode	377
TertiaryPhenotype	367
PrimaryPhenotype	364
ObservedFeatures	350
Glucose_mgdl	322
BP_systolic	321
Height_cm	320
BP_diastolic	318
VisitDate	316
Weight_kg	315
Creatinine_mgdl	307
Cholesterol_mgdl	299
SampleCollectionDateTime	295

Records Flagged for Missing Data: 2485

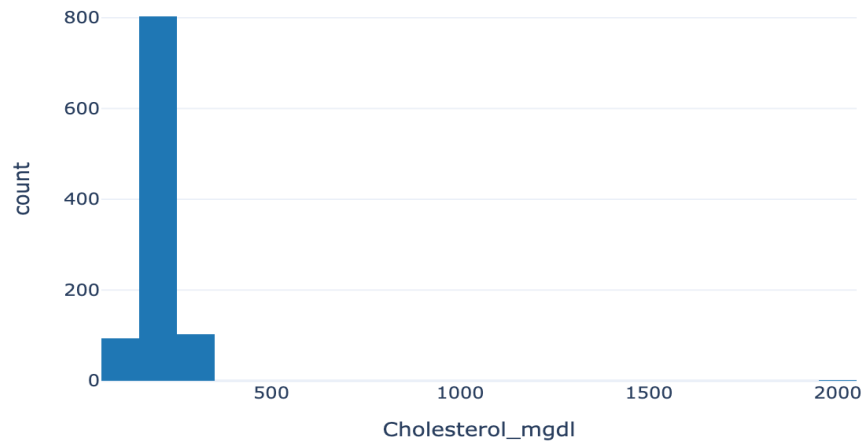
Ontology Mapping Summary

Ontology	Total Terms	Mapped	Success Rate
HPO	1985	1980	99.75%

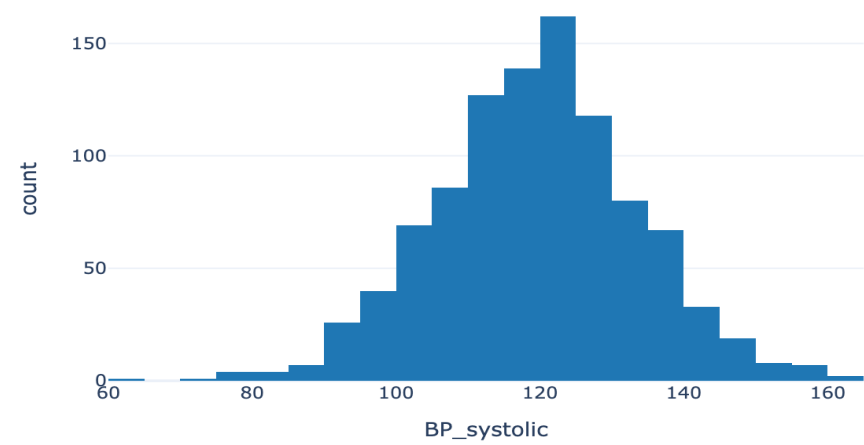
Visualizations



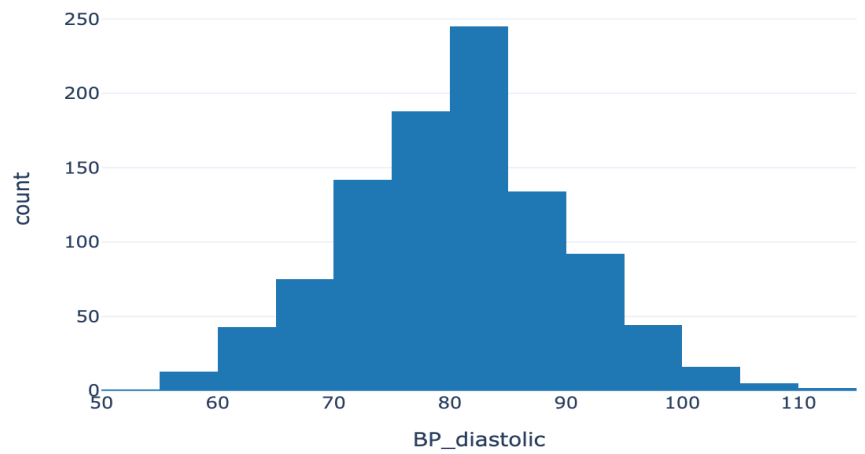
Distribution of Cholesterol_mgdl



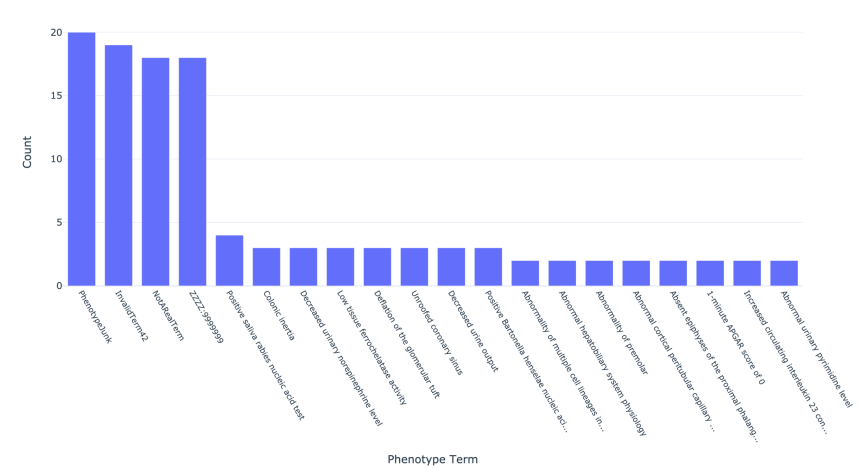
Distribution of BP_systolic



Distribution of BP_diastolic



Top 20 Most Common Terms in PrimaryPhenotype



Mapping Results: PrimaryPhenotype → HPO

