

Dataset Analysis Report

Dataset Path: binary_class_dataset

Dataset Overview:

- Total samples: 7,908
- Total entities: 308,853
- Number of splits: 3
- Unique labels: 3

Split Details:

- train: 6,326 samples
- validation: 790 samples
- test: 792 samples

Label Information:

- B-sensitive_data: 40,302 occurrences
- I-sensitive_data: 41,472 occurrences
- O: 227,079 occurrences

Split Statistics

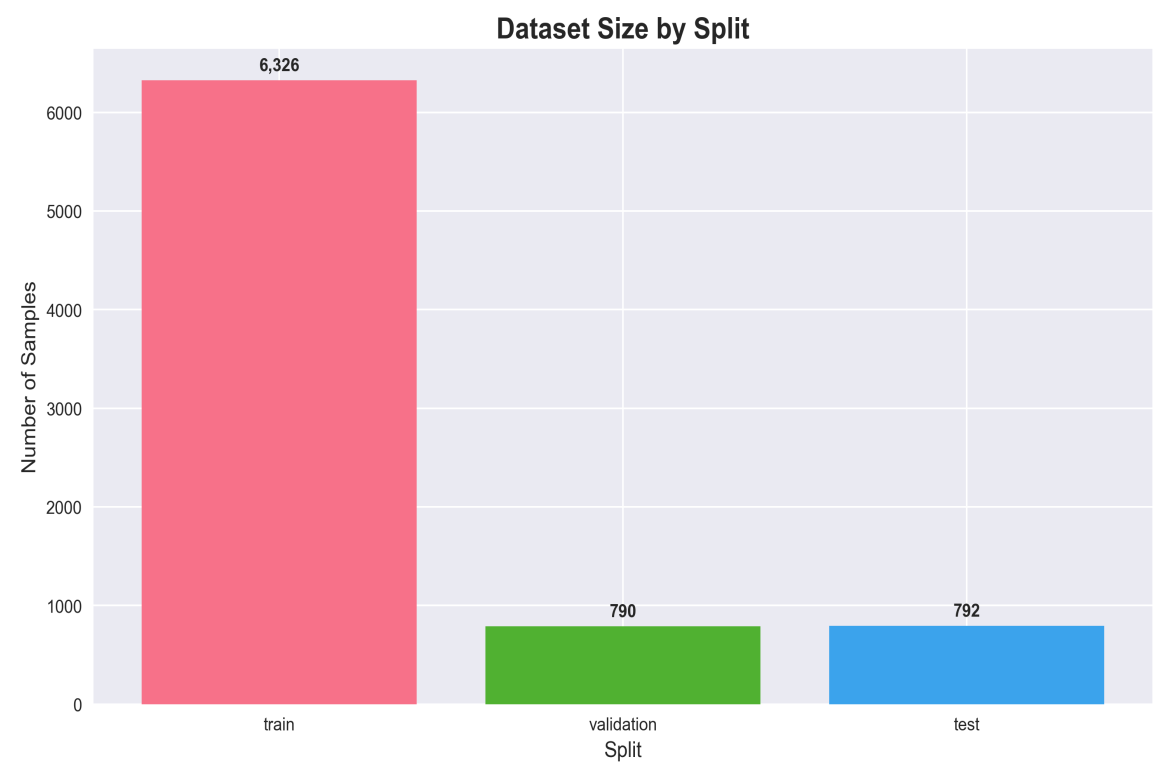
Split	Samples	Entities
train	6,326	247,118
validation	790	31,105
test	792	30,630

Label Distribution by Split

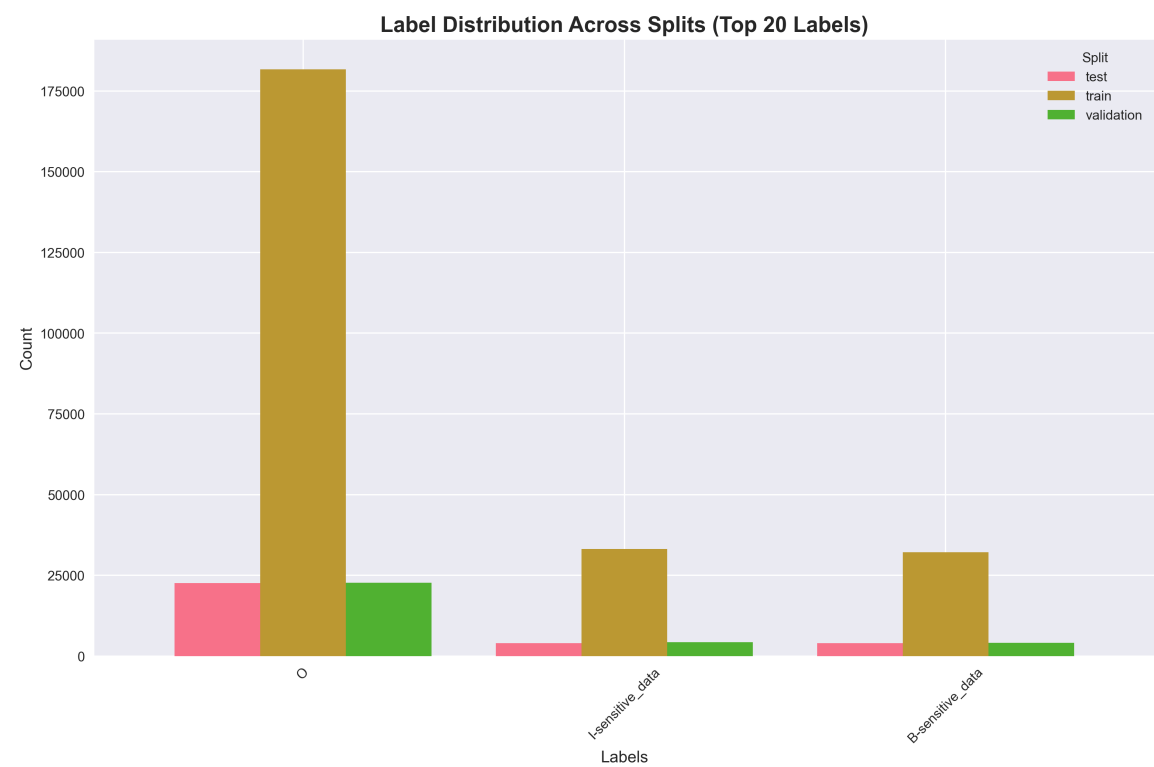
Label	train	validation	test
B-sensitive_data	32,200	4,131	3,971
I-sensitive_data	33,180	4,260	4,032
O	181,738	22,714	22,627

Visualizations

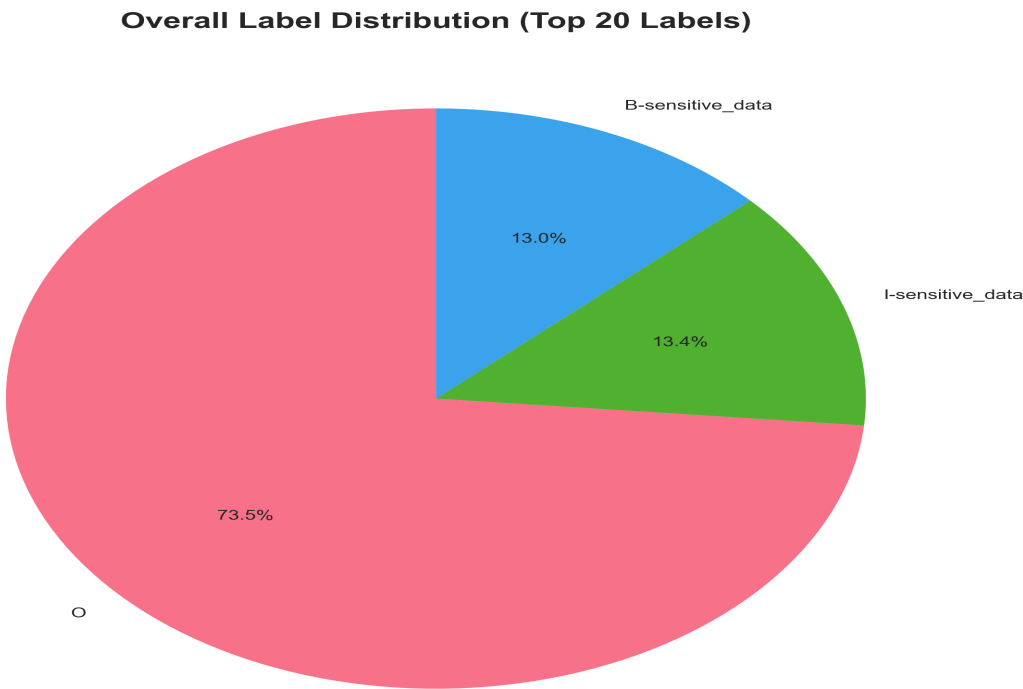
Dataset Size Comparison by Split



Label Distribution Across Splits (Top 20 Labels)



Overall Label Distribution (Top 20 Labels)



Total Label Frequencies (Top 20 Labels)

