

## 1 Datasets

Here we describe the datasets employed in our experiments.

### 1.1 *Hospitalization and Severity*

Both datasets were extracted from the same data source containing the symptoms and comorbidities declared when a citizen took a COVID test in the city of São José de Campos located in São Paulo state in Brazil. This data source was obtained through a partnership with the health Secretariat of São José dos Campos. The dataset encompasses the tests between March of 2020 and May of 2021. This database contains sensitive information and for ethical reasons can not be shared.

The *Hospitalization* dataset has three balanced classes: short (1:5 days), medium (6:10 days) and long (greater than 10 days) hospitalization. Table 2 summarize dataset information according to each class.

**Table 1.** Distribution of attributes per class in *Hospitalization* data set.

	SHORT	MEDIUM	LONG
NUMBER OF INSTANCES	1062	812	954
FEVER	0.573	0.596	0.591
COUGH	0.748	0.736	0.768
SORE THROAT	0.247	0.225	0.205
DYSPNEA	0.707	0.778	0.752
RESPIRATORY DISTRESS	0.529	0.499	0.534
LOW OXYGEN SATURATION	0.569	0.638	0.666
DIARRHEA	0.134	0.109	0.113
VOMIT	0.065	0.578	0.055
OTHER SYMPTOMS	0.652	0.639	0.600
CHRONIC CARDIOVASCULAR DISEASE	0.034	0.298	0.0437
IMMUNODEFICIENCY IMMUNODEPRESSION	0.028	0.030	0.021
CHRONIC KIDNEY DISEASE	0.030	0.020	0.019
DIABETES MELLITUS	0.266	0.284	0.334
ASTHMA	0.034	0.026	0.035
CHRONIC NEUROLOGICAL DISEASE	0.024	0.039	0.026
OTHER CHRONIC PNEUMOPATHY	0.034	0.030	0.034
OBESITY	0.059	0.070	0.073
OTHER RISKS	0.328	0.372	0.410
CHRONIC RESPIRATORY DISEASE	0.019	0.021	0.016
SEX (FEMALE)	0.455	0.432	0.413
AGE (MEAN - SD)	53.3 - 19.7	58.5 - 16.5	61.3 - 14.7

Dataset *Severity* has two classes: severe and non-severe patients. Severe patients are those who progressed to death within days 30 of hospitalization or

were hospitalized for ten or more days. Non severe patients were not hospitalized or were hospitalized for less than 10 days. This dataset was very imbalanced, with only 0.9% of instances in the severe class. Since the class imbalance affects the instance hardness analysis, in this work we randomly removed part of the majority class before the experiments. Table 2 summarizes dataset information according to each class.

## 1.2 *Hospital 1 and Hospital 2*

Both datasets were extracted from the same data source containing the laboratorial test results during the COVID pandemic for patients admitted at two private hospitals in the city of São Paulo. This data source is open and publicly available at <sup>1</sup>. The preprocessed datasets are available in the Supplementary material. In these datasets all patients are hospitalized and severity is assigned to patients who progress to death (within a period of 30 days from hospitalization) or for which the hospitalization period was larger than 14 days. Both datasets were imbalanced, and passed had instances randomly removed before the experiments. These datasets also contain missing values, which are input by the mean of the three nearest neighbors values. Table 4 summarizes some datasets information according to each class, after the imputation step.

---

<sup>1</sup> <https://repositoriodatasharingfapesp.uspdigital.usp.br/>

**Table 2.** Distribution of attributes per class in *Severity* dataset.

	SEVERE	NON-SEVERE
NUMBER OF INSTANCES	1710	1710
FEVER	0.551	0.680
COUGH	0.735	0.791
SORE THROAT	0.184	0.546
DYSPNEA	0.784	0.308
RESPIRATORY DISTRESS	0.568	0.041
LOW OXYGEN SATURATION	0.692	0.046
DIARRHEA	0.091	0.012
VOMIT	0.056	0.007
OTHER SYMPTOMS	0.556	0.750
CHRONIC CARDIOVASCULAR DISEASE	0.470	0.150
IMMUNODEFICIENCY IMMUNODEPRESSION	0.030	0.008
DIABETES MELLITUS	0.367	0.115
OBESITY	0.074	0.009
OTHER RISKS	0.438	0.029
CHRONIC RESPIRATORY DISEASE	0.017	0.036
SEX (FEMALE)	0.414	0.536
AGE (MEAN - SD)	64.7 - 15.7	42.0 - 16.0

**Table 3.** Summary of the datasets *Hospital 1* including number o patients, period of data collection, percentage of male patients, age statistics and blood tests.

	SEVERE	NON-SEVERE
PERIOD	03/2020 - 05/2021	
NUMBER OF INSTANCES	529	529
SEX (FEMALE)	0.357	0.340
AGE (MEAN - SD)	68.9 - 12.9	59.4 - 15.4
BLOOD TESTS	SODIUM, POTASSIUM, C-REACTIVE PROTEIN, GOT, GPT, UREA, D-DIMER, TROPONIN, CREATINE KINASE, BLOOD COUNT, CREATININE	

**Table 4.** Summary of the datasets *Hospital 2* including number o patients, period of data collection, percentage of male patients, age statistics and blood tests.

	SEVERE	NON-SEVERE
PERIOD	03/2020 - 01/2021	
NUMBER OF INSTANCES	67	67
SEX (FEMALE)	0.328	0.537
AGE (MEAN - SD)	65.8 - 14.6	55.2 - 17.0
BLOOD TESTS	SODIUM, POTASSIUM, C-REACTIVE PROTEIN, GOT, GPT, UREA, D-DIMER, TROPONIN, CREATINE KINASE, BLOOD COUNT, CREATININE	