

# Sarcopenia Dataset

## Writting Report: Statistical Analysis

Hanna Gloyna

### 1 Statistical Analysis

The used data set contained initially 250 entries, where each entry represents a patient, with 84 features. After removing ID, 17 entries were deleted as they were duplicates. Furthermore 8 misses in *Marcha* (engl. walking speed) were imputed with the mean. The encoding of the categorical data was changed, so that afterwards 0 did always correspond to either *Not answered*, *Not required* or *Regular value*. To obtain comparable values for variance of each feature a min-max-normalisation was applied, given by Equation 1.

$$\hat{x} = \frac{x - \min(x)}{\max(x) - \min(x)} \quad (1)$$

Where  $x$  corresponds to the original value of a feature,  $\min(x)$  the minimal values of this feature and  $\max(x)$  to the maximal value, respectively. Afterwards all features are within  $[0, 1]$ .

In accordance with specialists every representation of MMSE, Barthel, Norton, Lawton, and MNA, except the binary, was removed to reduce the dimensionality of the dataset. Additionally MM was dropped, because of its high correlation to IMM.

We dropped further features with too low variance or entropy, i.e. 90% of the values were the same. Additionally 7 features were deleted, because they contained a lot of misses or were included in another feature. Afterwards again all duplicates were dropped.

The remaining 231 entries included 42 serve and 189 mild sarcopenia cases. The dataset consists of 166 women and 65 men between 60 to 97

Feature 1	Feature 2	$r_s$
OA	Visiòn	-0.769
MED8	MED7	0.759
Drogas	Visiòn	-0.703

Table 1: Features with high correlation coefficient

Feature	p	f1
TNCM	0.000	14.971
Demencia	0.001	15.254
FuerzaPrension	0.006	8.957
Pantorrilla	0.097	40.598
IMM	0.135	24.095
LawtonAR	0.171	3.151
Dolor	0.173	2.867
Congiciòn	0.174	4.143
Charlson	0.178	2.372
CorreccionVisual	0.195	9.427

Table 2: Features and their corresponding p-value and f1-value to *sarcopenia*

years old and a median age of 79. By definition a serve sarcopenia cases is a decrease of IMM, grip strength and walking pace. If only a decrease of IMM and grip strength or walking speed is detected, it is classified as a mild case. With this definition 24 of 42 serve cases can be selected

The features with the highest Spearman correlation coefficient  $r_s$  are listed in Table 1. The used threshold for a high correlation was  $|r_s| > 0,7$  as suggested by Akoglu in [1]. In each row of Table 1 either feature 1 or feature 2 could be removed.

The results of a  $\chi^2$  and ANOVA analysis are displayed in Table 2. The 10 features with the low-

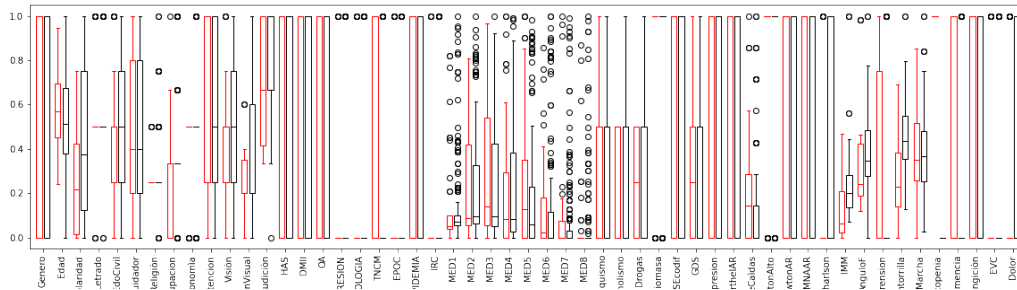


Figure 1: Boxplot of all remaining features, black boxes are created from all mild cases and red from all severe, respectively

est p-value from  $\chi_2$  analysis are almost the same as the ones with the highest f1 score. Only Occupation, NortonAlto and Edad (engl. age) achieved a higher f1 score, but a had also a higher p-value.

Figure 1 displays the plotted boxes of all remaining 53 features, where black boxes are generated from the subset of mild cases and red boxes from serve cases, respectively. Features with a difference at their mean and inter quartile range are especially interesting and very likely to be important for decision making. As can be seen features such as TNCM, Charlson, IMM, FuerzaPrension, Demencia and Dolor vary a lot for serve and mild cases. This also supports the calculated values and assignment of importance from the calculation of the p- and f1-value.

## 2 Hyperparameter Tuning

A *Support Vector Machine* (hereinafter: SVM) was used to predict the severity of sarcopenia. A cross validated grid search was applied to get an impression which parameters might be a good fit. Fig. 2 displays the results, where  $C$ ,  $\gamma$  and 3 different kernels were considered. In every setting the *Radial Basis Function* (rbf) did perform at least as good as the two other kernel. Furthermore it can be clearly seen that  $C = 2$  and  $\gamma = scale$  are the best setting with the rbf kernel.

### 3 Feature Selection

## References

- [1] Haldun Akoglu. “User’s guide to correlation coefficients”. In: *Turkish Journal of Emergency Medicine* 18.3 (Sept. 2018), pp. 91–93. ISSN: 24522473. DOI: 10.1016/j.tjem.2018.08.001. URL: <https://linkinghub.elsevier.com/retrieve/pii/S2452247318302164> (visited on 03/23/2022).

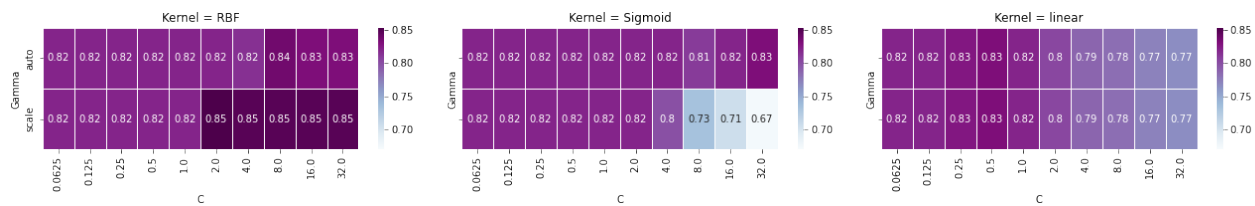


Figure 2: Results of hyperparameter tuning for SVM