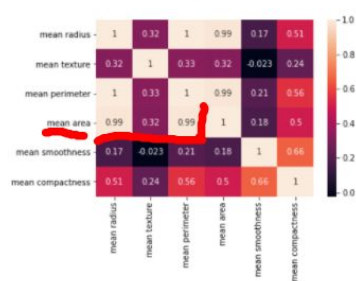


Estudo do efeito da normalização no cálculo da correlação

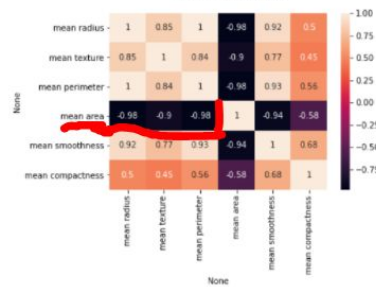
- Dataset com features sobre câncer de mama
- Dados em diferentes escalas
- Dados com distribuição assimétrica (skewed -> outliers)

Caso 1: sem processamento x
preprocessing.normalize x
preprocessing.scale

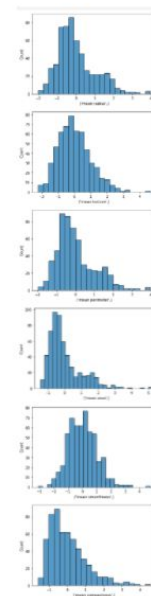
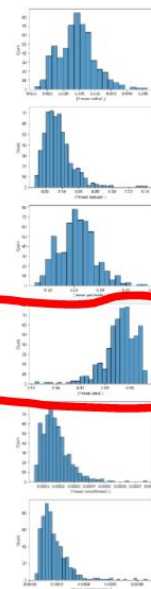
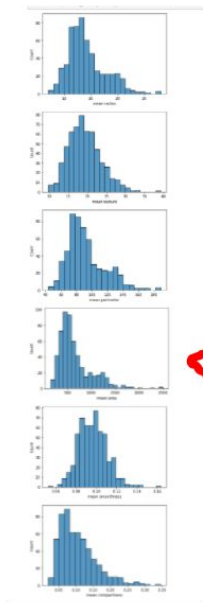
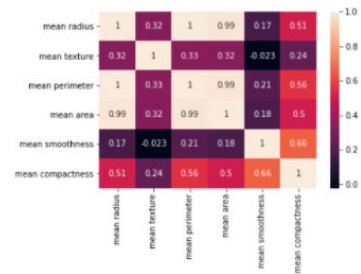
Sem proc.



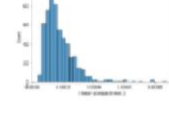
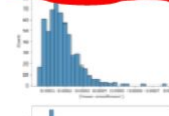
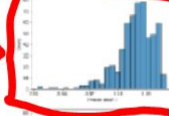
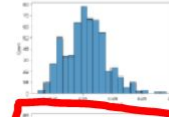
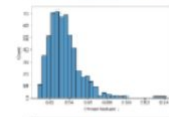
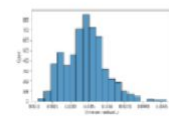
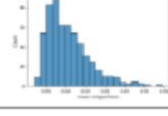
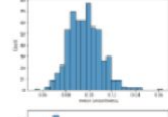
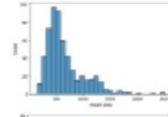
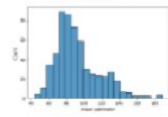
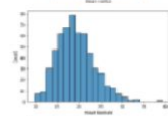
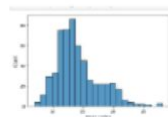
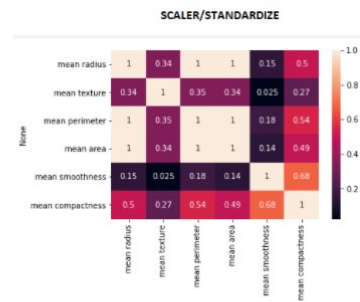
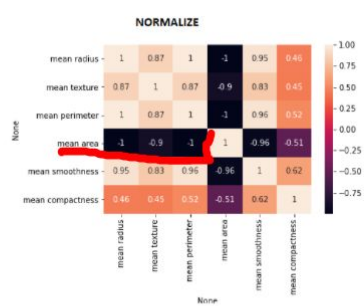
preprocessing.normalize



preprocessing.scale



PEARSON

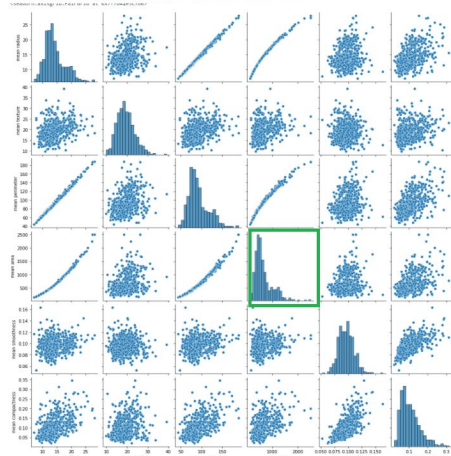


SPEARMAN

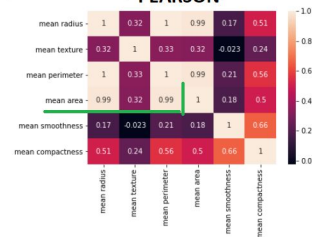
Caso 2: Sem processamento x

Transformação Logarítmica para normalização

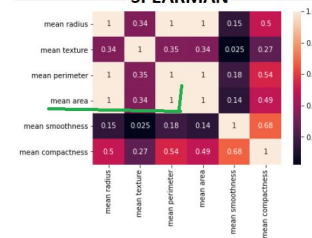
SEM PRE PROCESSAMENTO



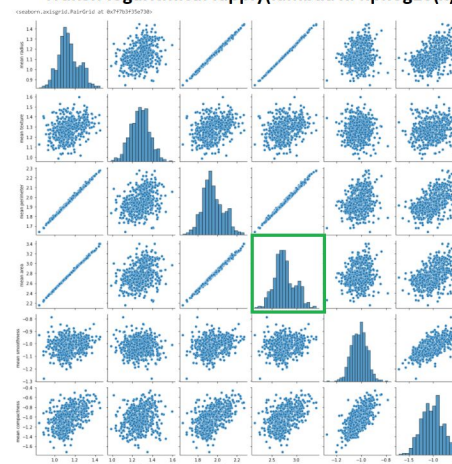
PEARSON



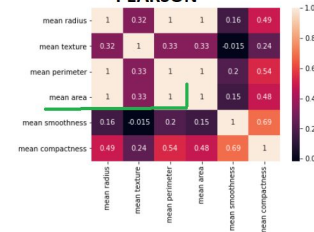
SPEARMAN



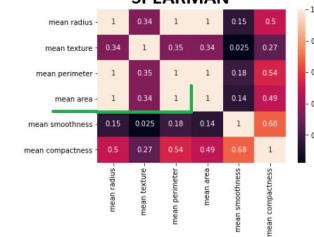
Transf. logarítmica: .apply(lambda x: np.log10(x))



PEARSON

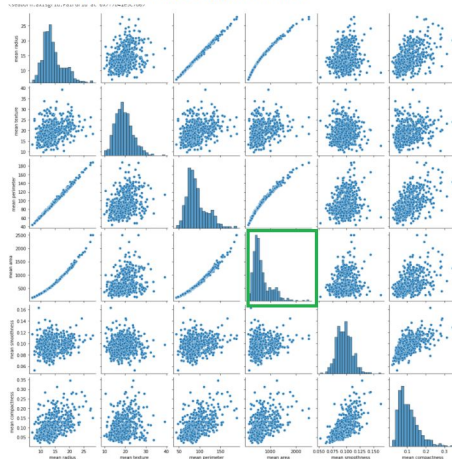


SPEARMAN

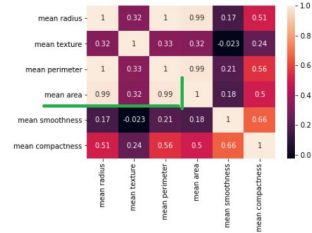


Caso 3: Sem processamento x
MinMax Scaler para normalização

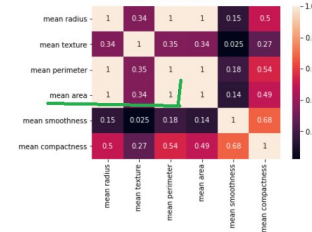
SEM PRE PROCESSAMENTO



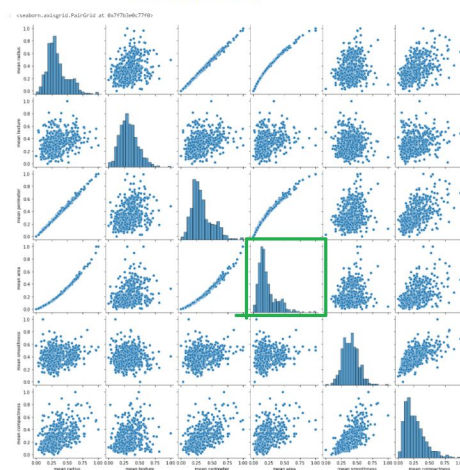
PEARSON



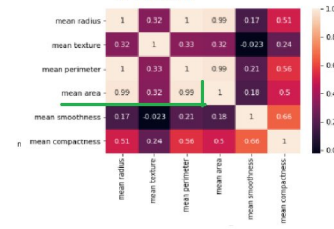
SPEARMAN



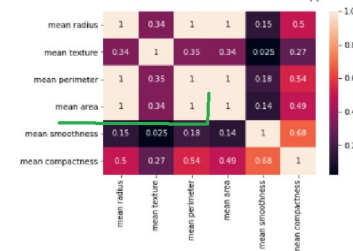
MinMax Scaler



PEARSON



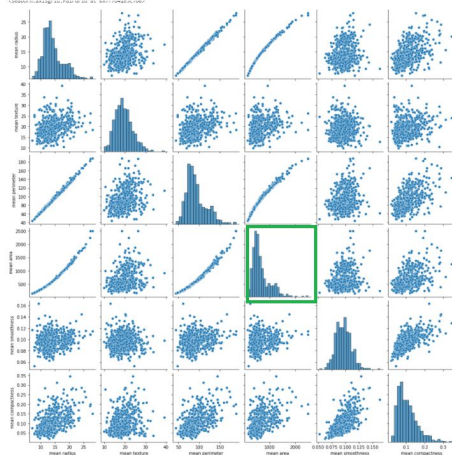
SPEARMAN



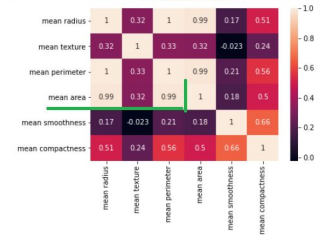
Caso 4: Sem processamento x Robust Scaler

- `with_centering = True, with_scaling = True`

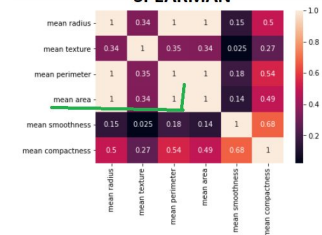
SEM PRE PROCESSAMENTO



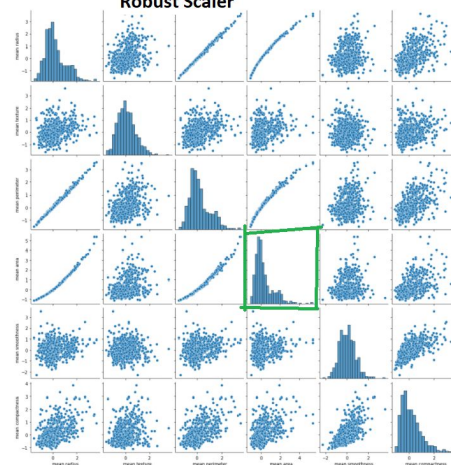
PEARSON



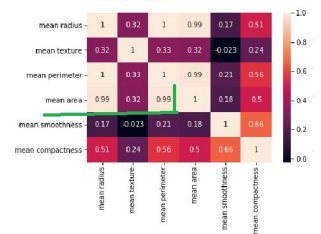
SPEARMAN



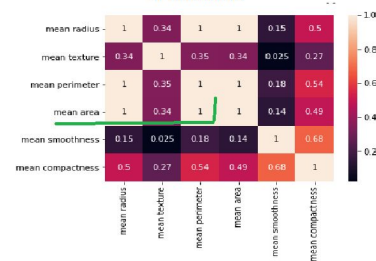
Robust Scaler



PEARSON

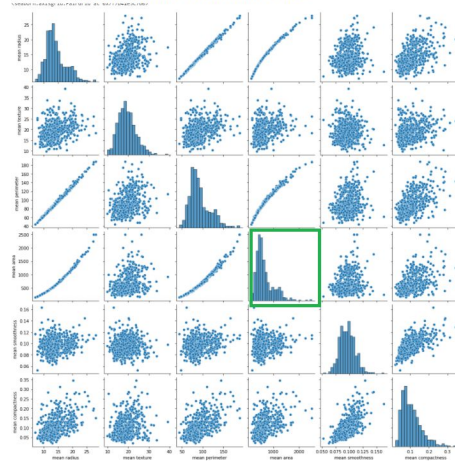


SPEARMAN

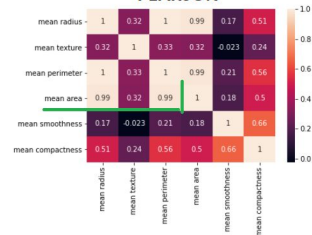


Caso 5: Sem processamento x Standard Scaler

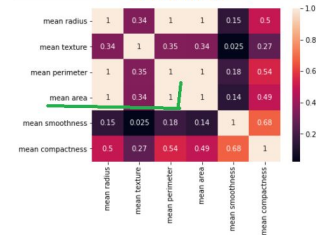
SEM PRE PROCESSAMENTO



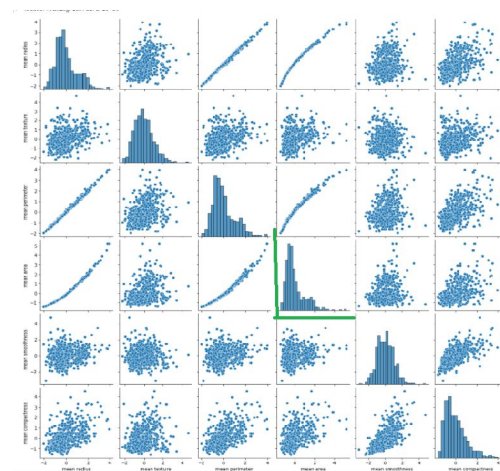
PEARSON



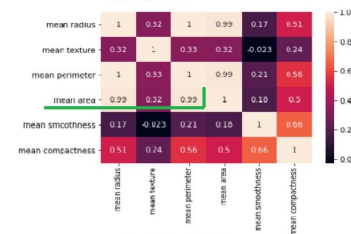
SPEARMAN



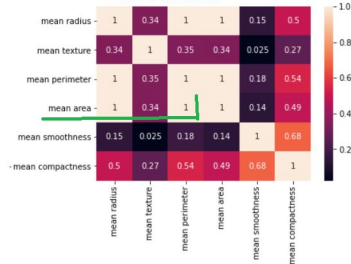
Standard Scaler



PEARSON



SPEARMAN



Conclusões

- Tanto normalização quanto escalonamento não parecem influenciar na correlação.
- A função `preprocessing.normalize` deve ser utilizada apenas para vetores :
<https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.normalize.html>
- Discussão sobre o assunto:

https://www.researchgate.net/post/Do_we_need_to_standardize_variables_with_different_scales_before_doing_correlation_analysis

$$\rho_{xy} = \text{Cov}(x, y) / (\sigma_x \sigma_y)$$

where:

ρ_{xy} = Pearson product-moment correlation coefficient

$\text{Cov}(x, y)$ = covariance of variables x and y

σ_x = standard deviation of x

σ_y = standard deviation of y

in this case, correlation is normalized by standard deviation. Therefore, no need to normalize them initially

[Cite](#)

Links utilizados

- https://www.researchgate.net/post/Do_we_need_to_standardize_variables_with_different_scales_before_doing_correlation_analysis
- <https://machinelearningmastery.com/standardscaler-and-minmaxscaler-transforms-in-python/>
- <https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.normalize.html>
- <https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.RobustScaler.html?highlight=robust%20scaler#sklearn.preprocessing.RobustScaler>
- https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.minmax_scale.html?highlight=min%20max#sklearn.preprocessing.minmax_scale