

Projetos em Ciência de Dados

Gianluca Devigili e Maisa O. Fraiz



1.

Preparação dos Dados

Tratamento de Dados

- + Dados salvos em **.pkl**
 - + Redução do tempo de carga de **28.4s** para **5.34s**
 - + Redução do uso de memória RAM
 - + Redução do espaço em disco de **3.8GB** para **3.4GB**

Higienização e Separação dos dados

- + Padronização dos nomes de colunas
- + Remoção de colunas que representam o mesmo dado (e.g. nome e sigla do time)
- + Redução do tamanho dos arquivos
- + Encapsulamento dos dados

Higienização e Separação dos dados



awards.pkl



games.pkl



playerBoxScores.pkl



players.pkl



playerTwitterFollowers.pkl



rosters.pkl



standings.pkl



targets.pkl



teamBoxScores.pkl



teams.pkl



teamTwitterFollowers.pkl



transactions.pkl

Higienização e Separação dos dados



awards.pkl



games.pkl



playerBoxScores.pkl



players.pkl



playerTwitterFollowers.pkl



rosters.pkl



standings.pkl



targets.pkl



teamBoxScores.pkl



teams.pkl

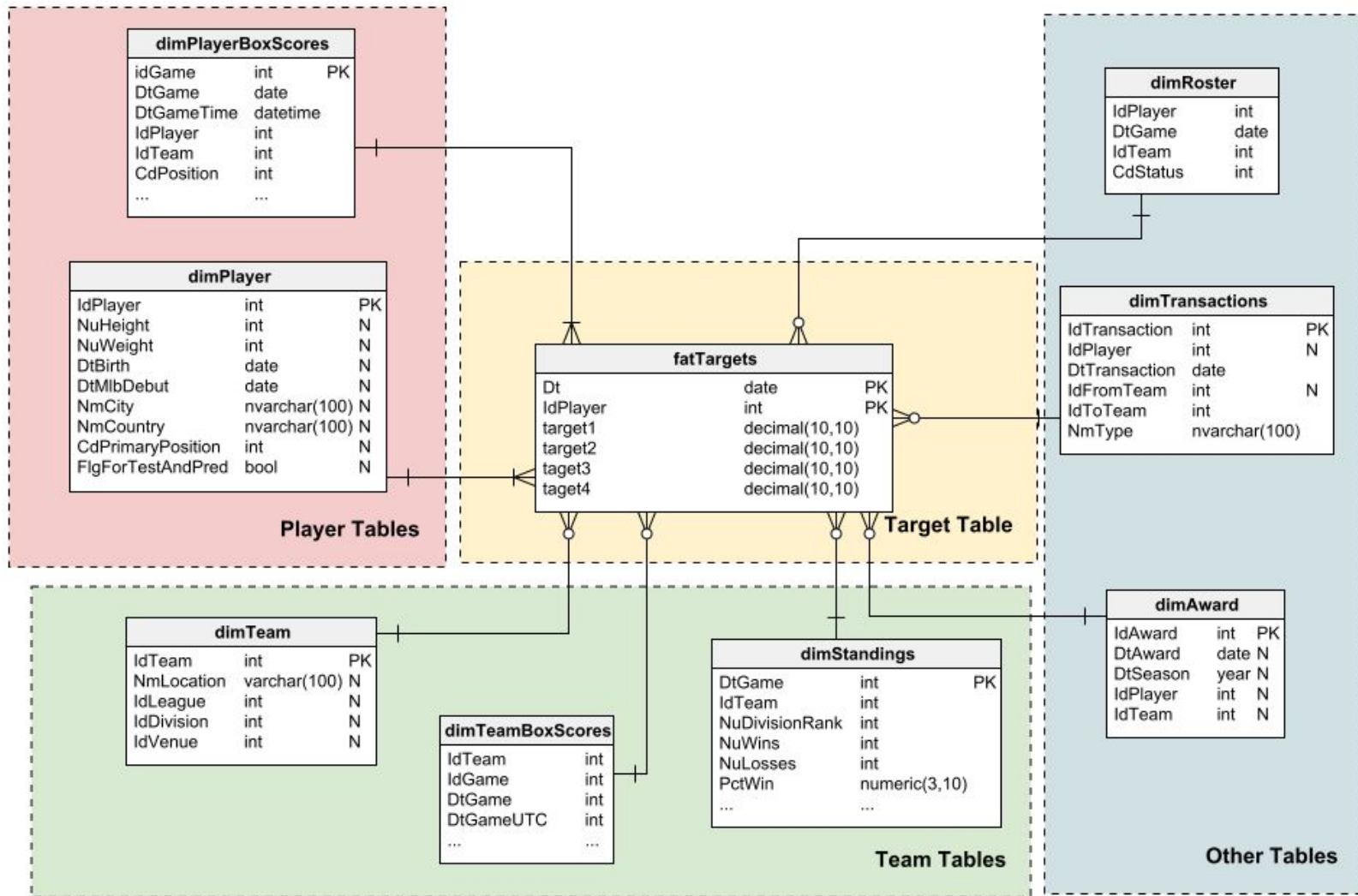


teamTwitterFollowers.pkl



transactions.pkl

Modelo do Data Warehouse





2.

Análise Exploratória

Estatísticas de Resumo dos Targets

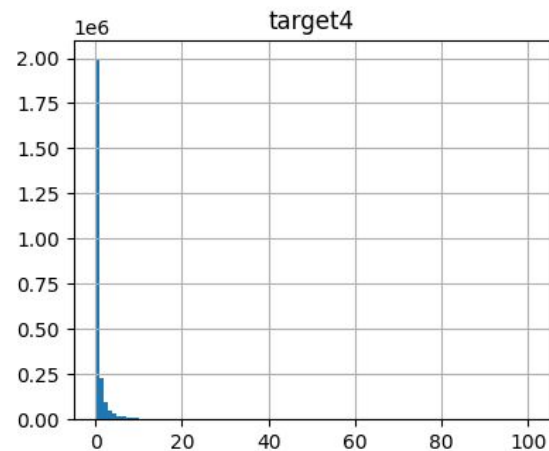
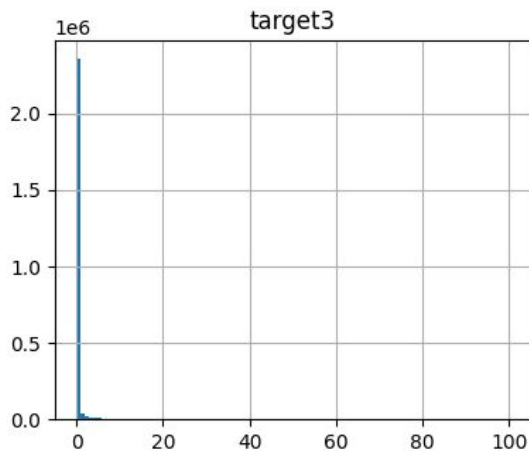
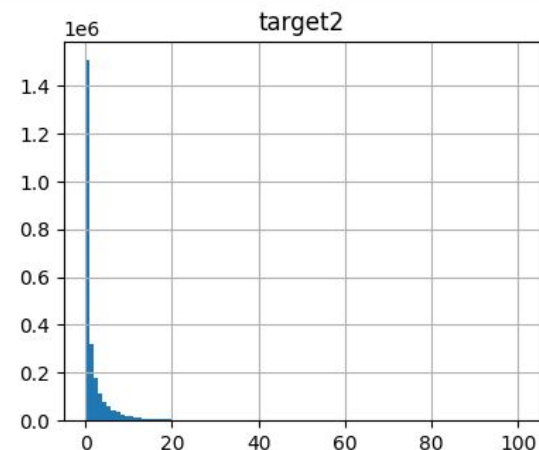
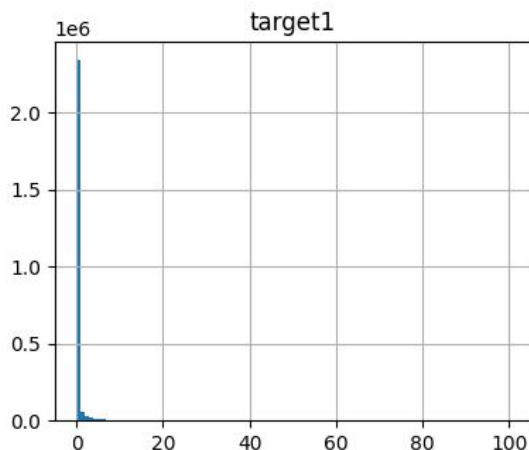
```
# describe the targets
trg_cols = ['target1', 'target2', 'target3', 'target4']
round(targets[trg_cols].describe(), 2)
```

	target1	target2	target3	target4
count	2506176.00	2506176.00	2506176.00	2506176.00
mean	0.57	2.46	0.69	1.14
std	4.17	6.23	5.07	4.23
min	0.00	0.00	0.00	0.00
25%	0.00	0.08	0.00	0.05
50%	0.00	0.56	0.00	0.22
75%	0.02	2.24	0.02	0.76
max	100.00	100.00	100.00	100.00

- + 2 milhões de registros
- + Os targets variam de 0 a 100
- + Médias baixas
- + Maioria dos valores 0 ou próximos de 0

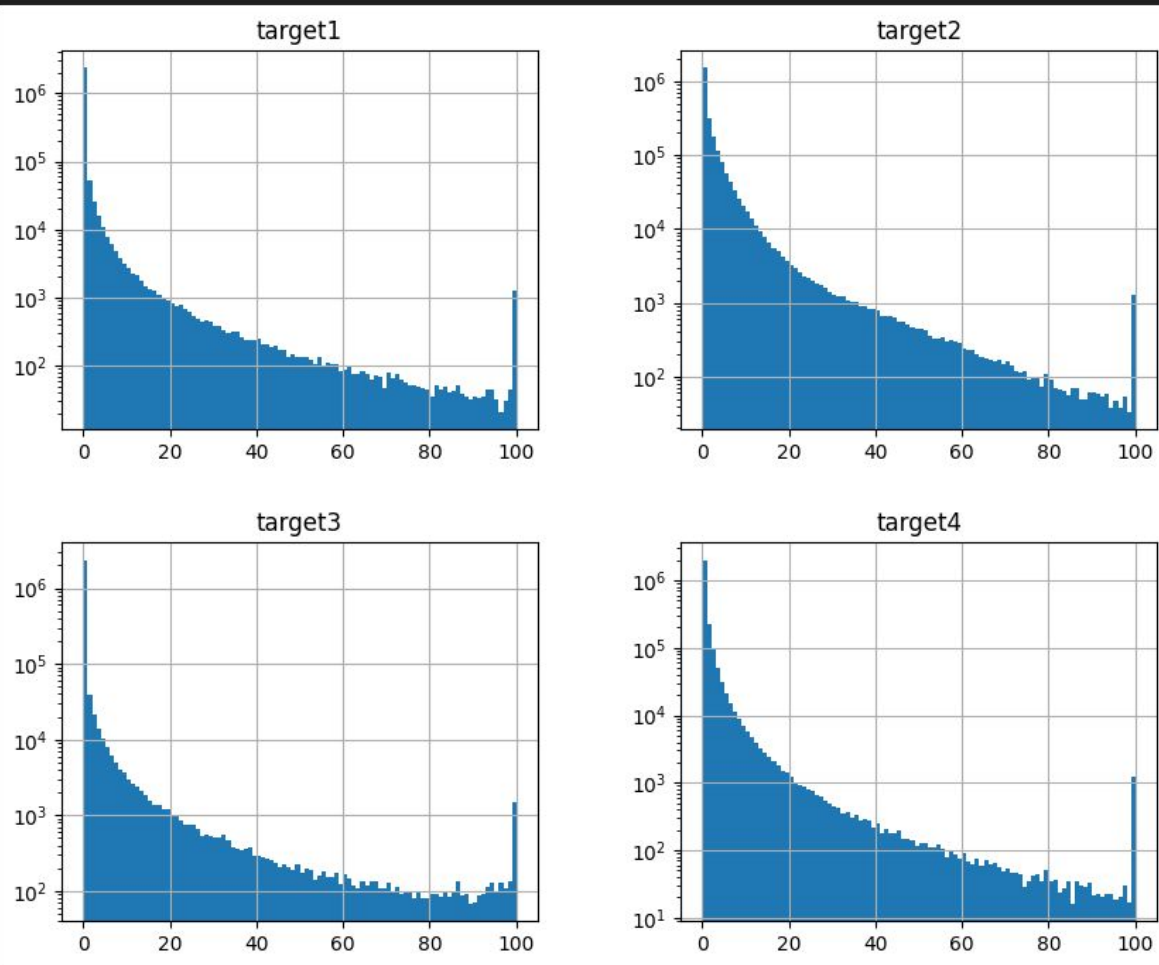
Histograma dos targets

```
targets[trg_cols].hist(figsize=(10, 8), bins=100);
```



Histograma dos targets em escala logarítmica

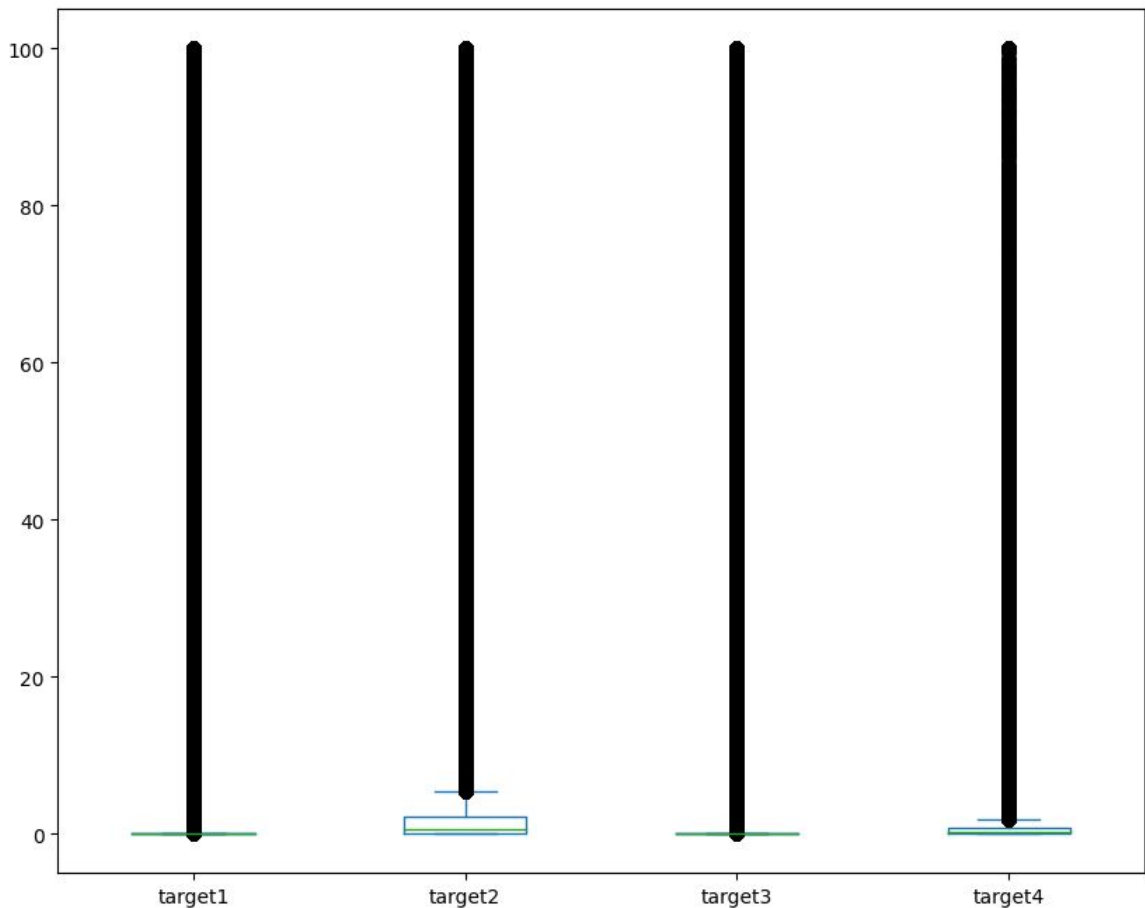
```
targets[trg_cols].hist(figsize=(10, 8), bins=100, log=True);
```



Boxplot dos Targets

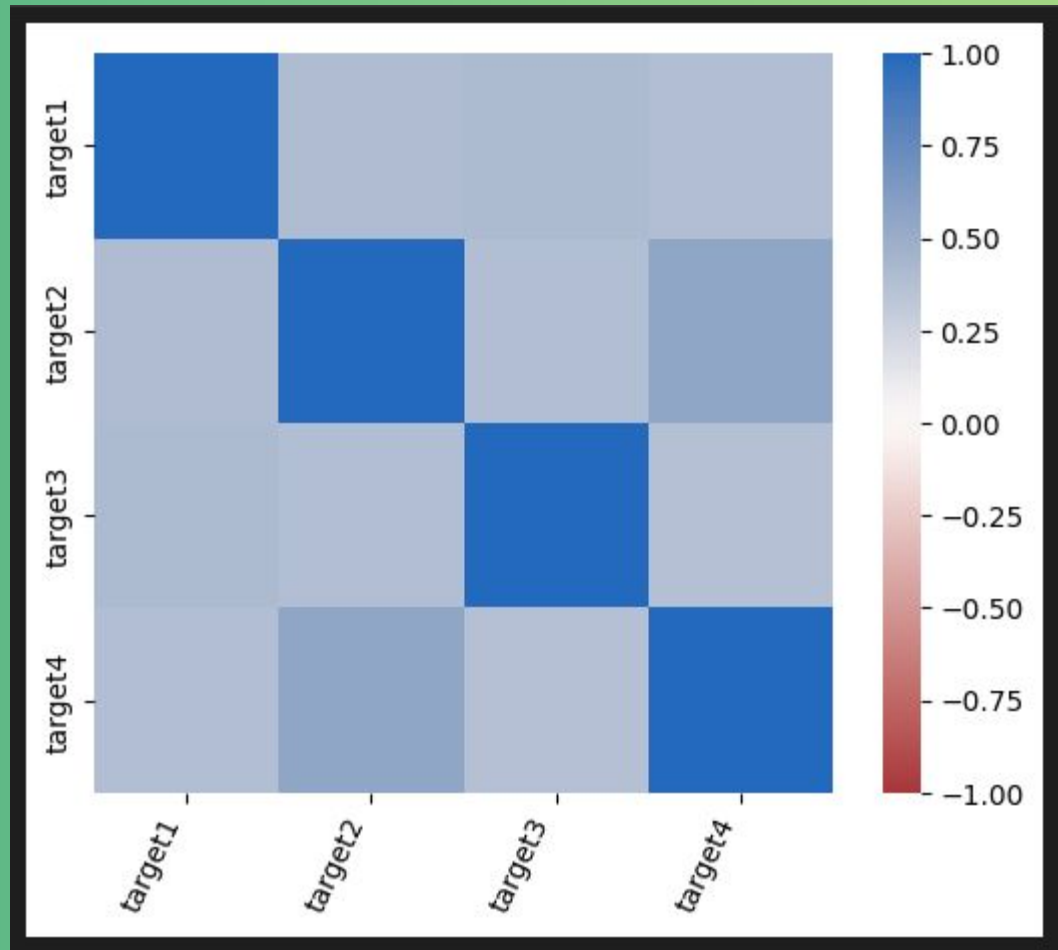
12

```
targets[trg_cols].plot(kind='box', figsize=(10, 8));
```

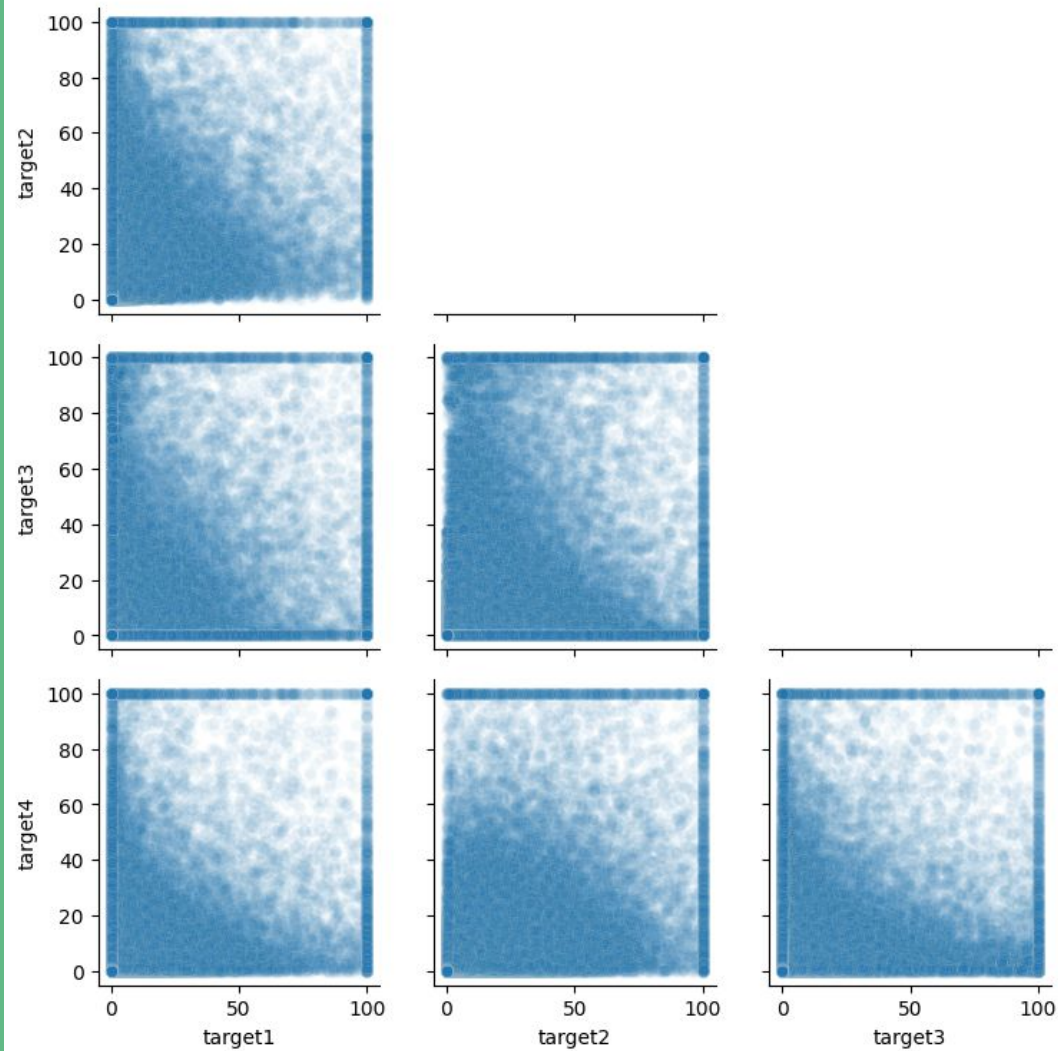


Correlação Entre os Targets

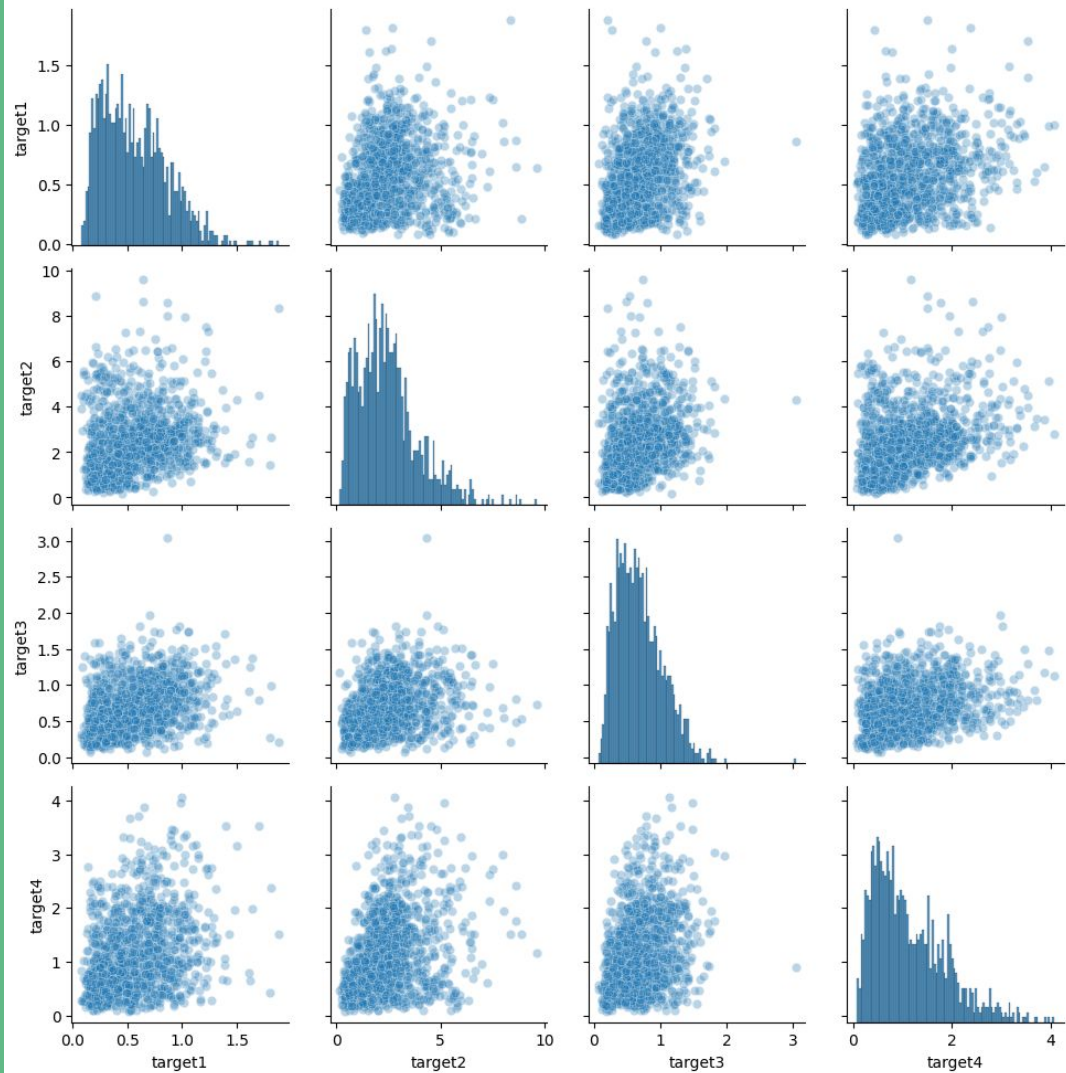
	target1	target2	target3	target4
target1	1.000000	0.404532	0.411024	0.384962
target2	0.404532	1.000000	0.388134	0.548991
target3	0.411024	0.388134	1.000000	0.370333
target4	0.384962	0.548991	0.370333	1.000000



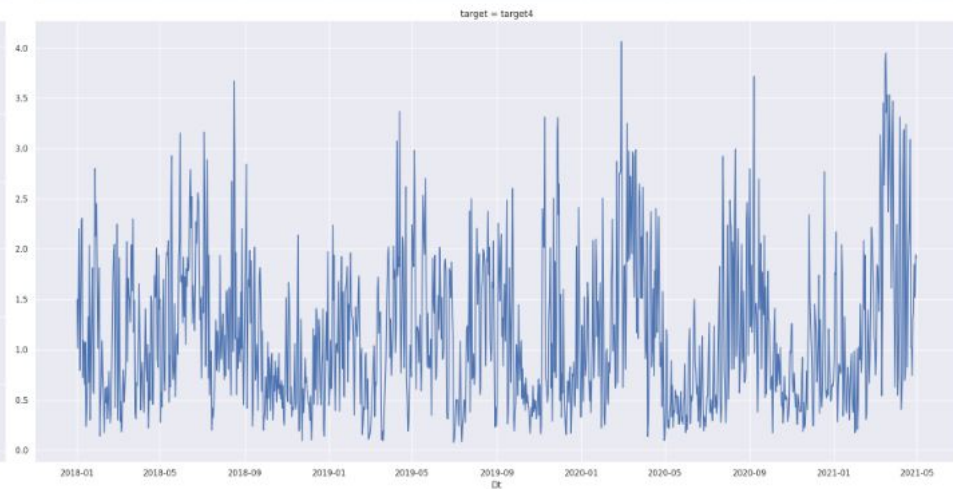
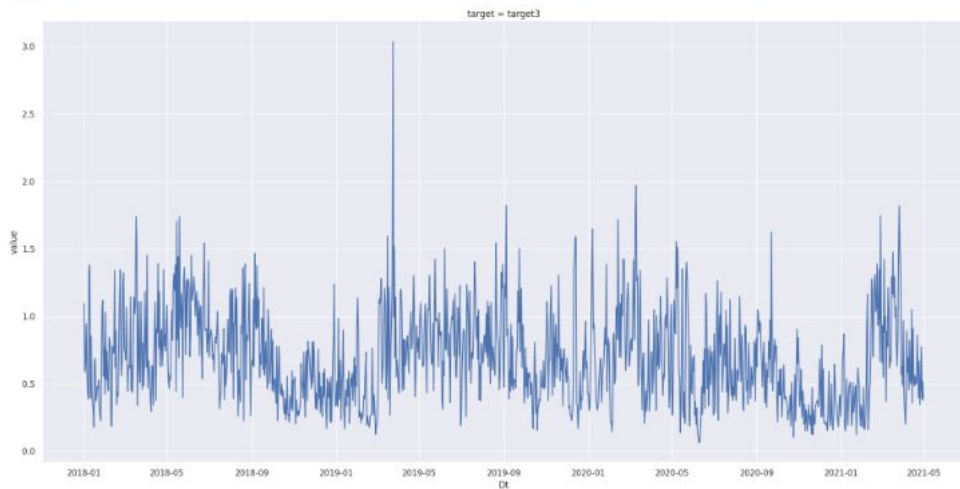
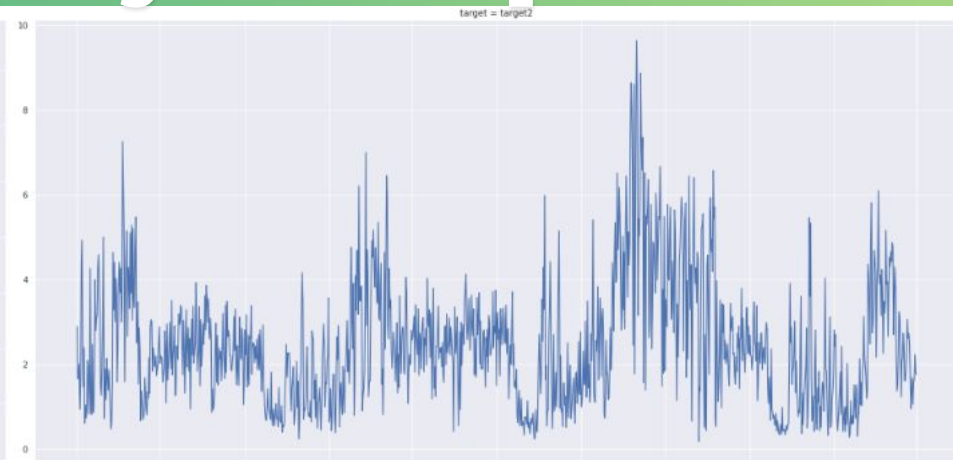
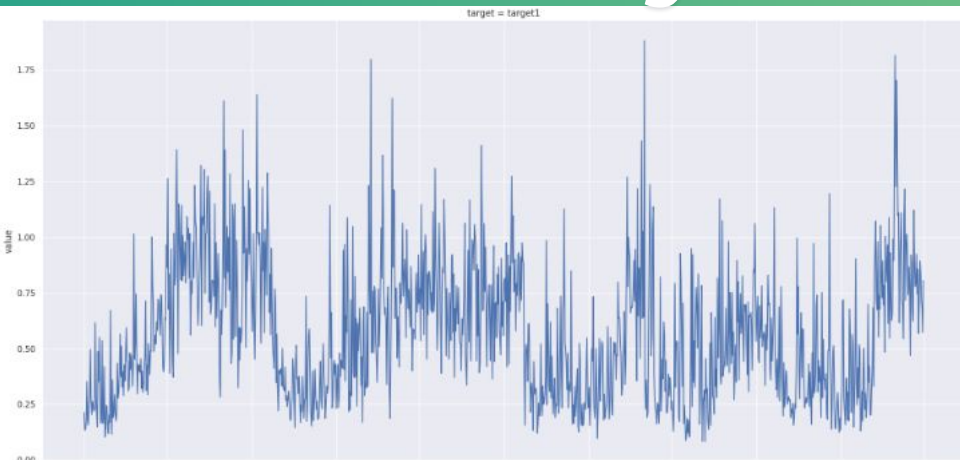
Scatterplot dos Targets



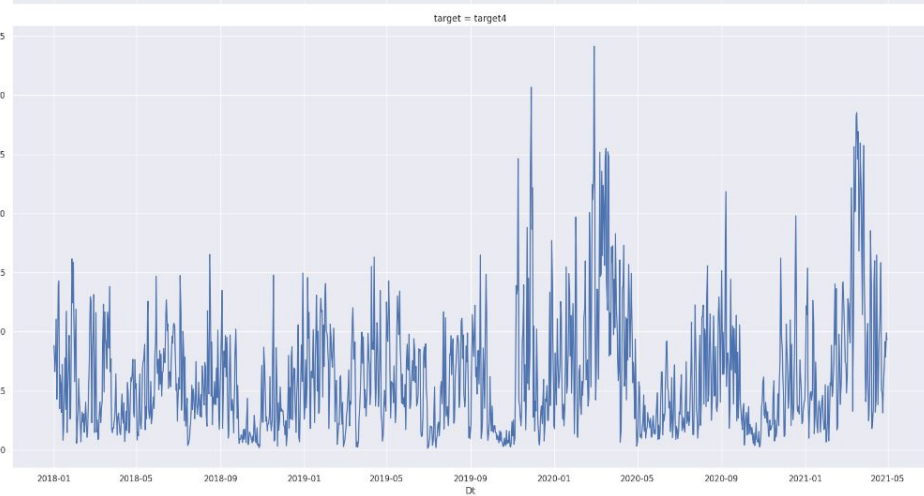
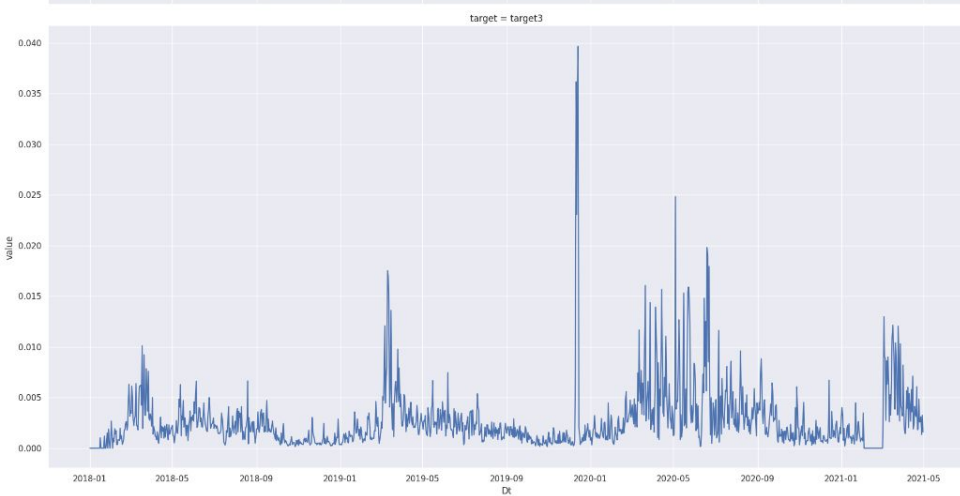
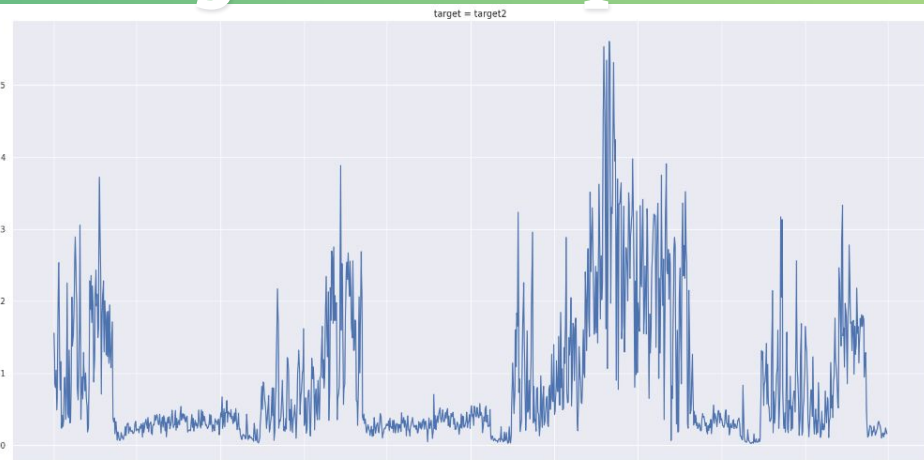
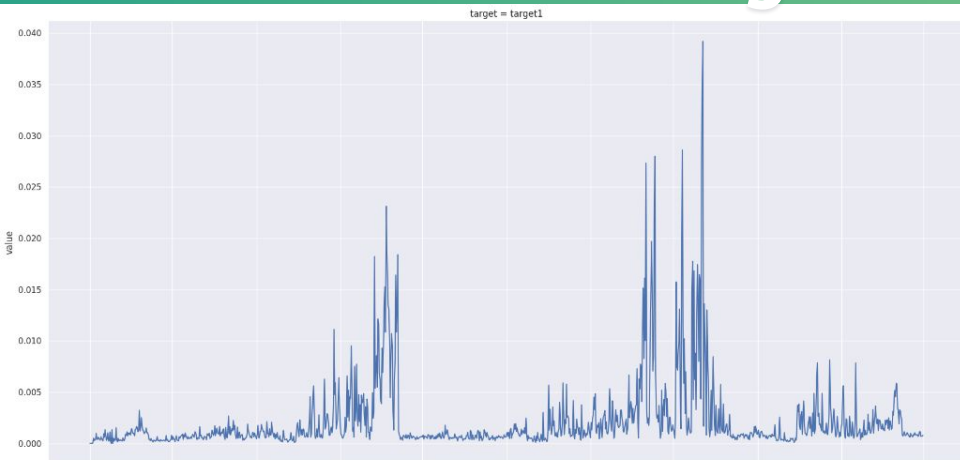
Média diária dos Targets

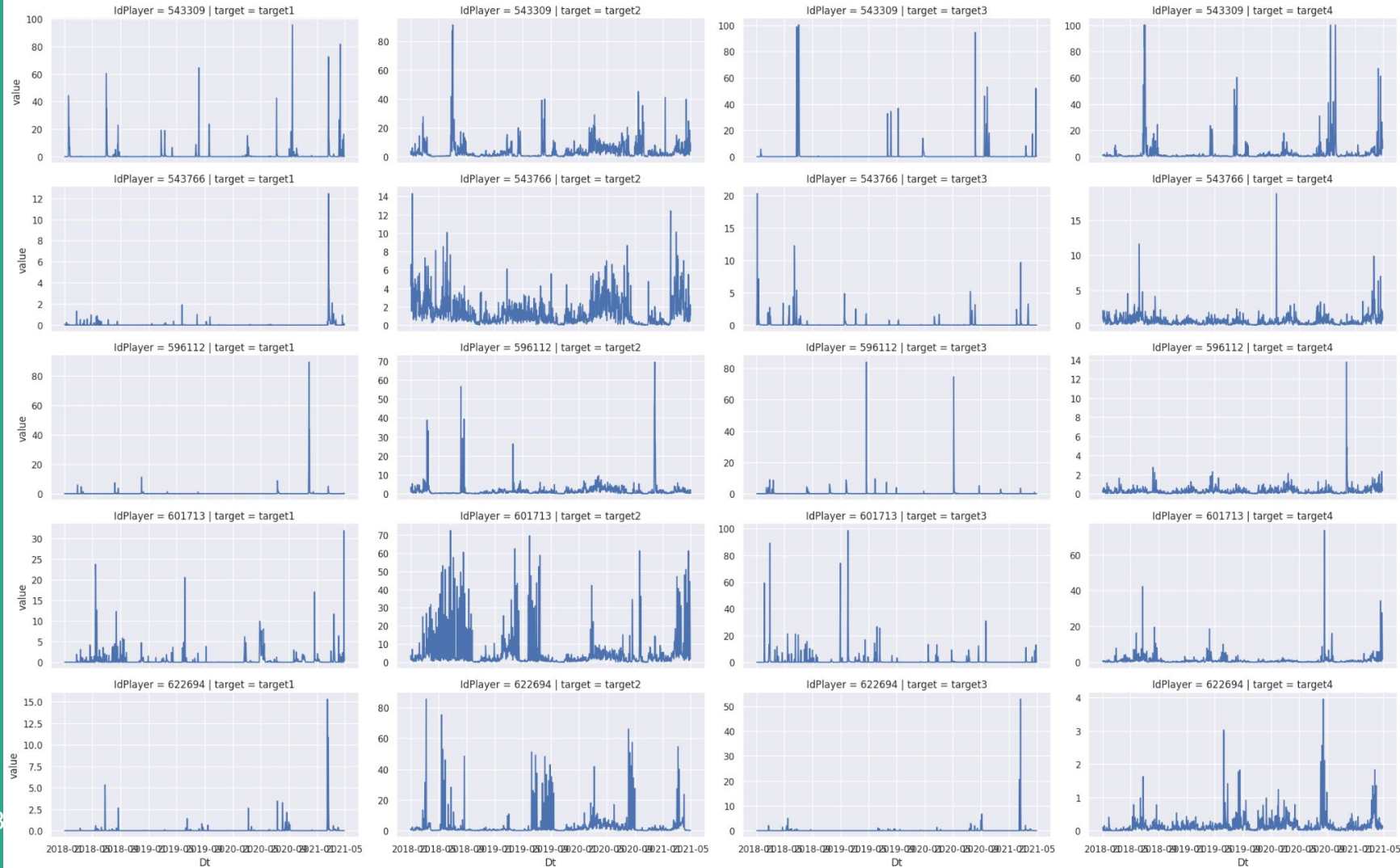


Média dos Targets ao longo do tempo



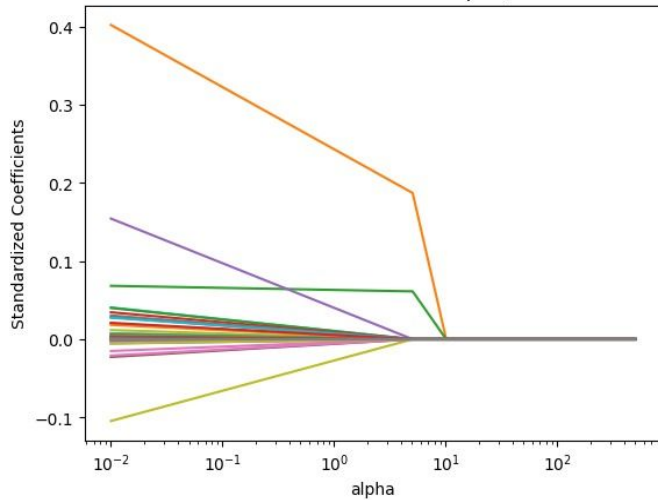
Mediana dos Targets ao longo do tempo



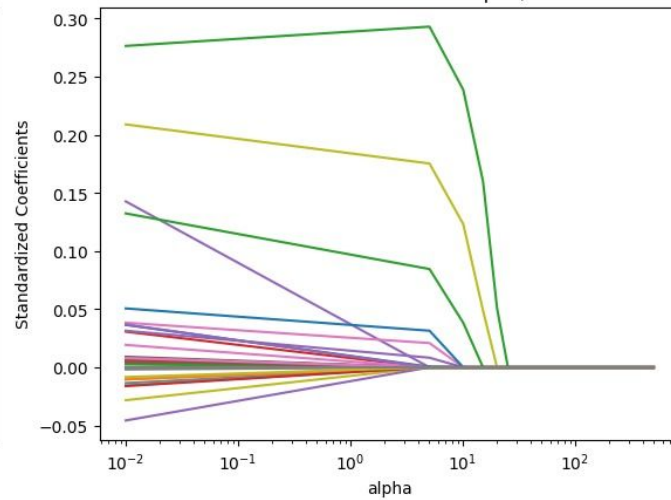




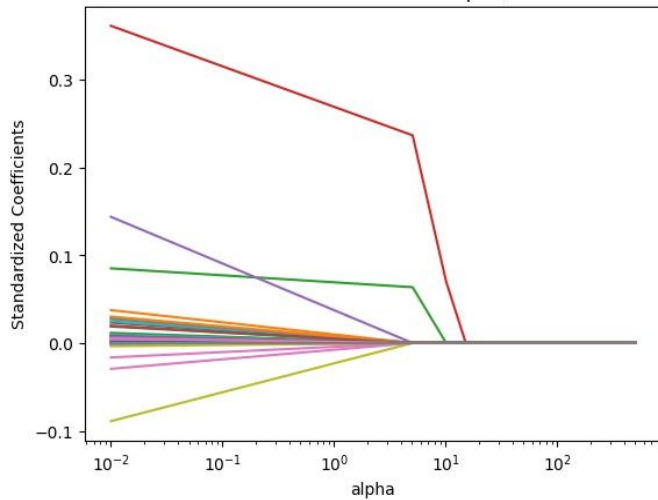
Lasso coefficients as a function of alpha, TARGET = 1



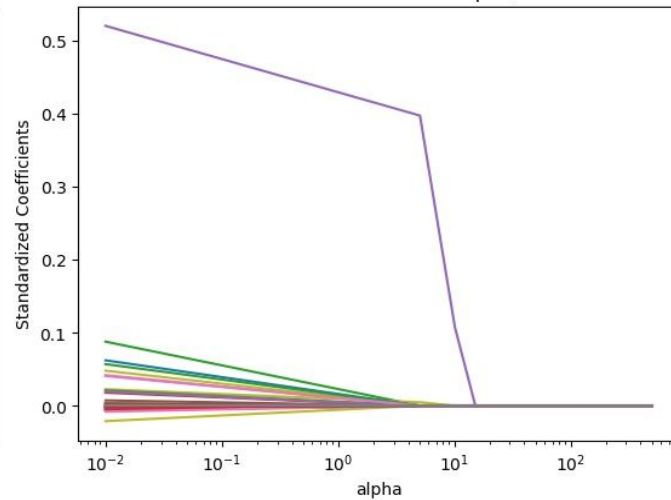
Lasso coefficients as a function of alpha, TARGET = 2



Lasso coefficients as a function of alpha, TARGET = 3



Lasso coefficients as a function of alpha, TARGET = 4



Projetos em Ciência de Dados

Gianluca Devigili e Maisa O. Fraiz