

```

1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 fonte = "https://github.com/alura-cursos/imersao-dados-2-2020/blob/master/MICRO
6
7 dados = pd.read_csv(fonte)

1 provas = ["NU_NOTA_CN", "NU_NOTA_CH", "NU_NOTA_MT", "NU_NOTA_LC", "NU_NOTA_REDAC
2 dados["NU_NOTA_TOTAL"] = dados[provas].sum(axis = 1)
3 provas = ["NU_NOTA_CN", "NU_NOTA_CH", "NU_NOTA_MT", "NU_NOTA_LC", "NU_NOTA_REDAC
4 dados_sem_notas_zero = dados.query("NU_NOTA_TOTAL != 0")
5 dados_sem_notas_NaN = dados_sem_notas_zero[provas].dropna()

1 provas_entrada = ["NU_NOTA_CH", "NU_NOTA_LC", "NU_NOTA_CN", "NU_NOTA_REDACAO"]
2 provas_saida = "NU_NOTA_MT"
3 notas_entrada = dados_sem_notas_NaN[provas_entrada]
4 notas_saida = dados_sem_notas_NaN[provas_saida]

1 x = notas_entrada
2 y = notas_saida

1 from sklearn.model_selection import train_test_split
2 SEED = 4321
3 x_treino, x_teste, y_treino, y_teste = train_test_split(x,y, test_size = 0.25, r

1 x_teste.shape

(23135, 4)

1 from sklearn.svm import LinearSVR
2
3 modelo = LinearSVR(random_state = SEED)
4 modelo.fit(x_treino, y_treino)

/usr/local/lib/python3.6/dist-packages/sklearn/svm/_base.py:947: ConvergenceWarning:
  "the number of iterations.", ConvergenceWarning)
LinearSVR(C=1.0, dual=True, epsilon=0.0, fit_intercept=True,
          intercept_scaling=1.0, loss='epsilon_insensitive', max_iter=1000,
          random_state=4321, tol=0.0001, verbose=0)

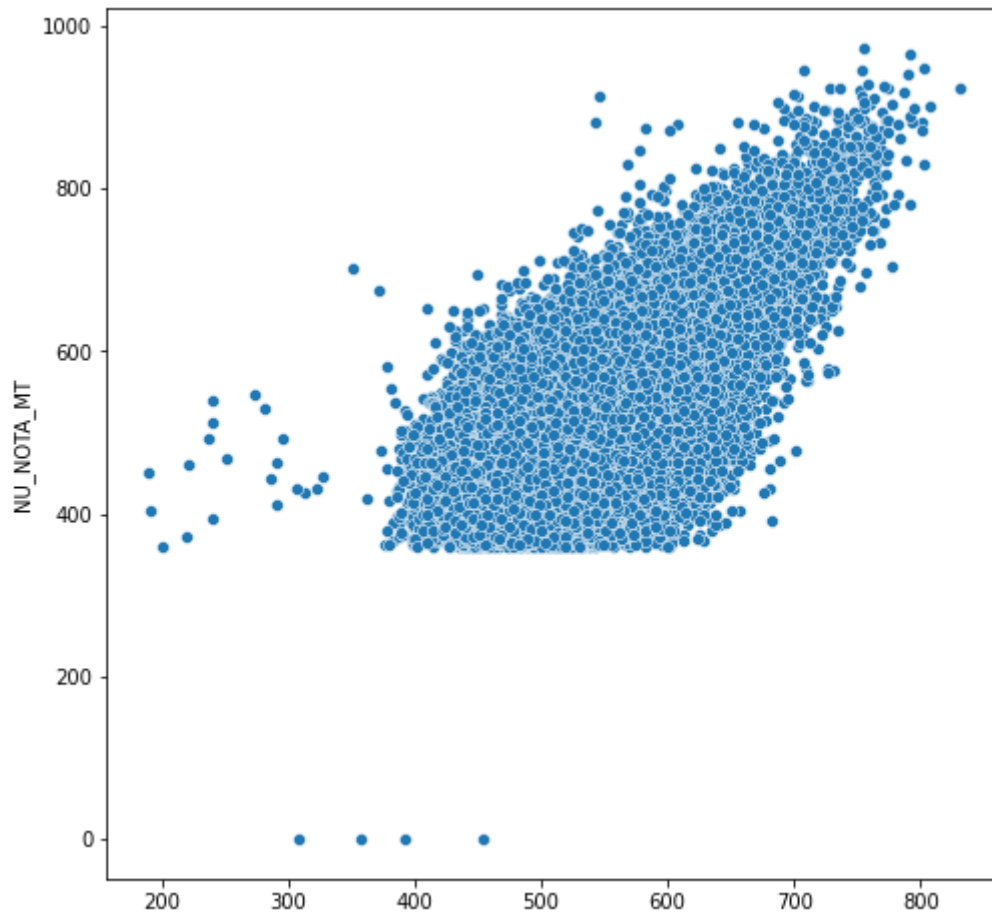
1 predicoes_MT = modelo.predict(x_teste)

1 y_teste[:5]
```

```
114991    459.7
104685    617.2
.....    ....
```

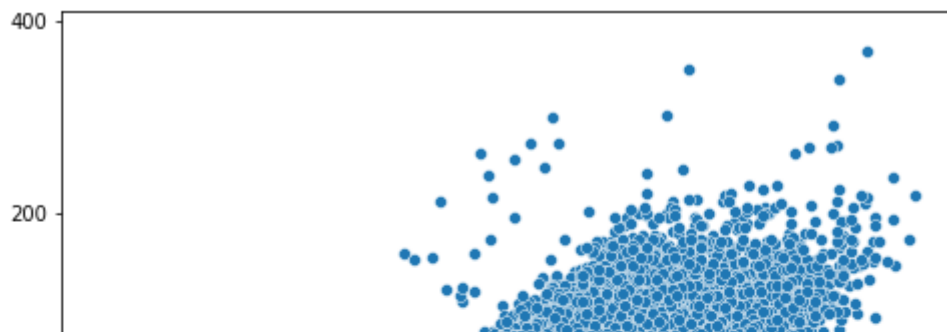
```
1 plt.figure(figsize=(8,8))
2 sns.scatterplot(x=predicoes_MT, y = y_teste)
3
4
```

<matplotlib.axes._subplots.AxesSubplot at 0x7f607bc41e10>



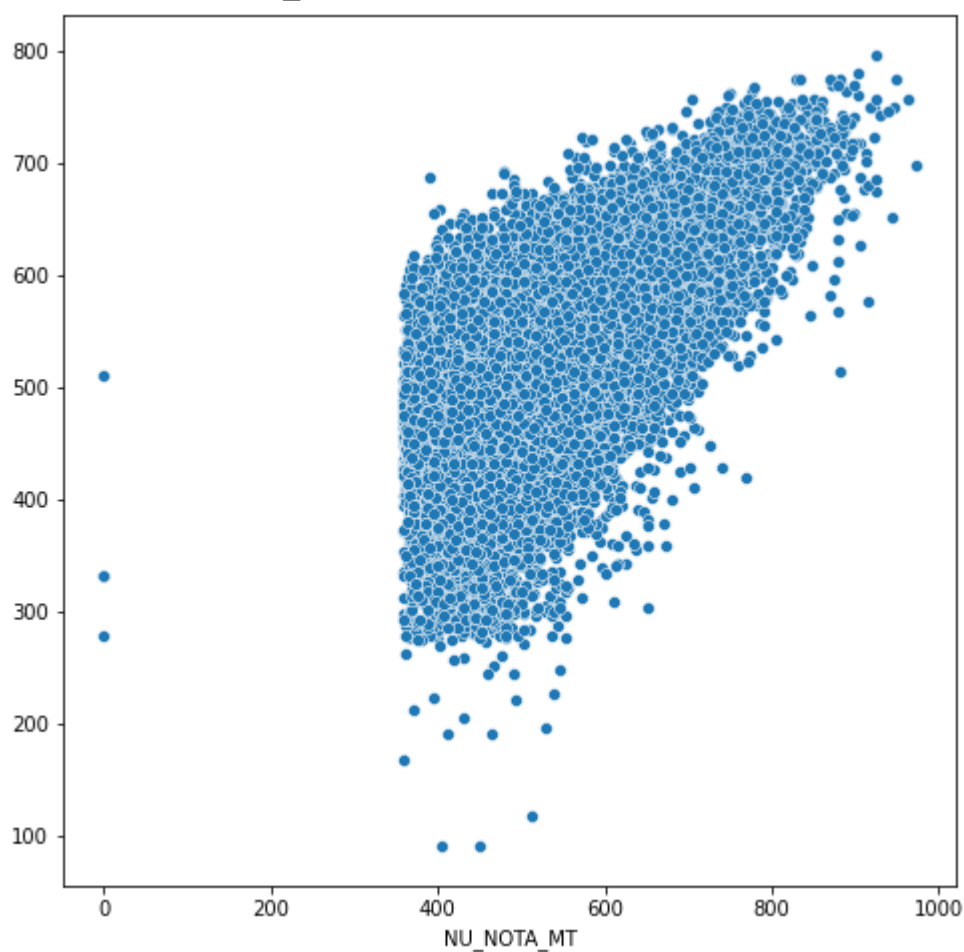
```
1 plt.figure(figsize=(8,8))
2 sns.scatterplot(x=y_teste, y = y_teste-predicoes_MT)
```

<matplotlib.axes._subplots.AxesSubplot at 0x7f60793a1cf8>



```
1 plt.figure(figsize=(8,8))
2 sns.scatterplot(x=y_teste, y = x_teste.mean(axis=1))
```

<matplotlib.axes._subplots.AxesSubplot at 0x7f60776d5ba8>



```
1 resultados = pd.DataFrame()
2 resultados["Real"] = y_teste
3 resultados["Previsao"] = predicoes_MT
4 resultados["diferenca"] = resultados["Real"] - resultados["Previsao"]
5 resultados["quadrado_diferenca"] = (resultados["Real"] - resultados["Previsao"])
6
```

```
1 resultados
```

	Real	Previsao	diferenca	quadrado_diferenca
114991	459.7	564.486105	-104.786105	10980.127878
104685	617.2	628.305027	-11.105027	123.321625
91028	520.0	509.980986	10.019014	100.380650
115802	703.8	593.681144	110.118856	12126.162420
93303	627.1	579.784588	47.315412	2238.748252
...
81393	688.6	626.291851	62.308149	3882.305380
40159	479.1	652.313922	-173.213922	30003.062777
64083	501.4	607.787172	-106.387172	11318.230438
84661	772.5	734.227663	38.272337	1464.771772

```
1 resultados["quadrado_diferenca"].mean()**(1/2)
```

```
81.26755383435341
```

```
1 from sklearn.dummy import DummyRegressor
```

```
2
```

```
3 modelo_dummy = DummyRegressor()
```

```
4 modelo_dummy.fit(x_treino, y_treino)
```

```
5 dummy_predicoes = modelo_dummy.predict(x_teste)
```

```
6
```

```
7 from sklearn.metrics import mean_squared_error
```

```
8
```

```
9 mean_squared_error(y_teste, dummy_predicoes) #12063.645588509502
```

```
12063.645588509502
```

```
1 mean_squared_error(y_teste, predicoes_MT)#6604.41530621957
```

```
6604.41530621957
```

Desafio1 = procurar outro modelo de ML para treinar e comparar com os modelos criados em aula

```
1 #desafio 1 com decision tree
```

```
2 from sklearn import tree
```

```
3 modelo_tree = tree.DecisionTreeRegressor()
```

```
4 modelo_tree.fit(x_treino, y_treino)
```

```
5 tree_predicoes = modelo_tree.predict(x_teste)
```

```
6 mean_squared_error(y_teste, tree_predicoes) #11365.011590231252
```

```
7
```

```
11382.14146401556
```

```

1 resultadosdt = pd.DataFrame()
2 resultadosdt["Real"] = y_teste
3 resultadosdt["Previsao"] = tree_predicoes
4 resultadosdt["diferenca"] = resultadosdt["Real"] - resultadosdt["Previsao"]
5 resultadosdt["quadrado_diferenca"] = (resultadosdt["Real"] - resultadosdt["Previsao"])**2
6 resultadosdt

```

	Real	Previsao	diferenca	quadrado_diferenca
114991	459.7	671.0	-211.3	44647.69
104685	617.2	664.6	-47.4	2246.76
91028	520.0	480.8	39.2	1536.64
115802	703.8	553.2	150.6	22680.36
93303	627.1	509.5	117.6	13829.76
...
81393	688.6	529.8	158.8	25217.44
40159	479.1	633.4	-154.3	23808.49
64083	501.4	660.0	-158.6	25153.96
84661	772.5	816.3	-43.8	1918.44
79106	369.0	558.7	-189.7	35986.09

23135 rows x 4 columns

```

1 #desafio 1 com near neighbors
2 from sklearn.neighbors import KNeighborsRegressor
3 modelo_neigh = KNeighborsRegressor()
4 modelo_neigh.fit(x_treino, y_treino)
5 neigh_predicoes = modelo_neigh.predict(x_teste)
6 mean_squared_error(y_teste, neigh_predicoes) #6432.6058905554355

```

6432.6058905554355

```

1 resultadosnn = pd.DataFrame()
2 resultadosnn["Real"] = y_teste
3 resultadosnn["Previsao"] = neigh_predicoes
4 resultadosnn["diferenca"] = resultadosnn["Real"] - resultadosnn["Previsao"]
5 resultadosnn["quadrado_diferenca"] = (resultadosnn["Real"] - resultadosnn["Previsao"])**2
6 resultadosnn

```

	Real	Previsao	diferenca	quadrado_diferenca
114991	459.7	465.70	-211.3	44647.69
104685	617.2	621.74	-47.4	2246.76
91028	520.0	495.72	39.2	1536.64
115802	703.8	506.04	150.6	22680.36
93303	627.1	521.90	117.6	13829.76
...
81393	688.6	554.86	158.8	25217.44

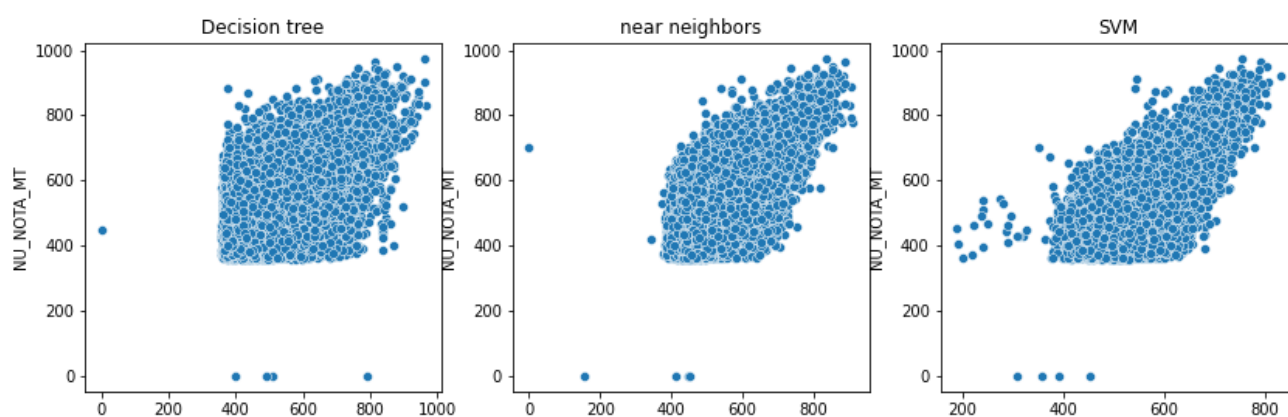
```

1 plt.figure(figsize=(14,14))
2 plt.subplot(3,3,1)
3 #arvore de decisoes
4 sns.scatterplot(x=tree_predicoes, y = y_teste)
5 plt.title("Decision tree")
6 #SVM
7 plt.subplot(3,3,3)
8 sns.scatterplot(x= predicoes_MT, y = y_teste)
9 plt.title("SVM")
10 #near neighbors
11 plt.subplot(3,3,2)
12 sns.scatterplot(x= neigh_predicoes, y = y_teste)
13 plt.title("near neighbors")
14
15 plt.suptitle("Comparando os modelos")

```

Text(0.5, 0.98, 'Comparando os modelos')

Comparando os modelos



Desafio2: Ler a documentacao do Dummy e alterar o metodo de regressao

```
1 #desafio 2, dummy strategy = mean
```

```
2 modelo_dummy_mean = DummyRegressor(strategy="mean")
3 modelo_dummy_mean.fit(x_treino, y_treino)
4 dummy_predicoes_mean = modelo_dummy_mean.predict(x_teste)
5 mean_squared_error(y_teste, dummy_predicoes_mean) #12063.645588509502
```

12063.645588509502

```
1 #desafio 2, dummy strategy =median
2 modelo_dummy_median = DummyRegressor(strategy="median")
3 modelo_dummy_median.fit(x_treino, y_treino)
4 dummy_predicoes_median = modelo_dummy_median.predict(x_teste)
5 mean_squared_error(y_teste, dummy_predicoes_median) #12599.969971795981
```

12599.969971795981

```
1 #desafio 2, dummy strategy =quantile - 1(maximo)
2 modelo_dummy_quantile = DummyRegressor(strategy="quantile", quantile = 1.0)
3 modelo_dummy_quantile.fit(x_treino, y_treino)
4 dummy_predicoes_quantile = modelo_dummy_quantile.predict(x_teste)
5 mean_squared_error(y_teste, dummy_predicoes_quantile ) #223890.15624378645
```

223890.15624378645

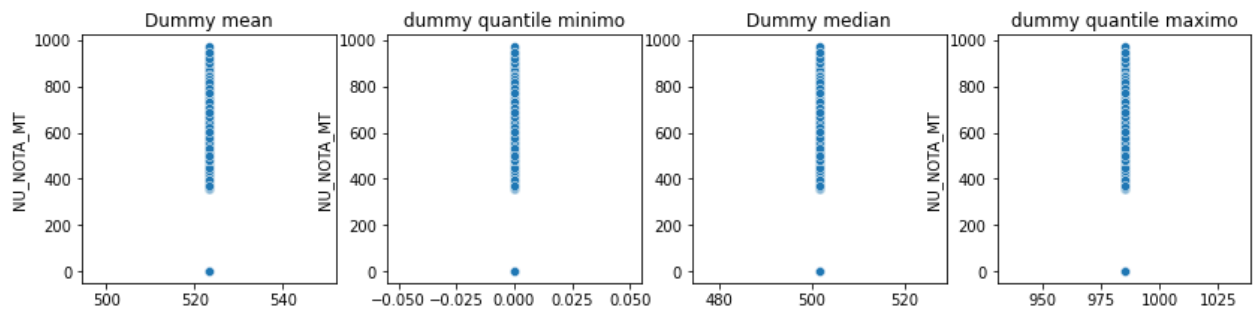
```
1 #desafio 2, dummy strategy =quantile - 0(minimo)
2 modelo_dummy_quantile_m = DummyRegressor(strategy="quantile", quantile = 0.0)
3 modelo_dummy_quantile_m.fit(x_treino, y_treino)
4 dummy_predicoes_quantile_m = modelo_dummy_quantile_m.predict(x_teste)
5 mean_squared_error(y_teste, dummy_predicoes_quantile_m ) #287425.98995893664
```

287425.98995893664

```
1 plt.figure(figsize=(14,14))
2 plt.subplot(4,4,1)
3 #mean
4 sns.scatterplot(x=dummy_predicoes_mean, y = y_teste)
5 plt.title("Dummy mean")
6 #median
7 plt.subplot(4,4,3)
8 sns.scatterplot(x= dummy_predicoes_median, y = y_teste)
9 plt.title("Dummy median")
10 #minimo
11 plt.subplot(4,4,2)
12 sns.scatterplot(x= dummy_predicoes_quantile_m, y = y_teste)
13 plt.title("dummy quantile minimo")
14
15 #maximo
16 plt.subplot(4,4,4)
17 sns.scatterplot(x= dummy_predicoes_quantile, y = y_teste)
18 plt.title("dummy quantile maximo")
19
20
21 plt.suptitle("Comparando os metodos")
```

Text(0.5, 0.98, 'Comparando os metodos')

Comparando os metodos



Desafio3: Busacar outra métrica(mean squared root) para avaliar modelos de regressão

```

1 import sklearn.metrics as skm
2 #Metricas para o SVM
3 n1 = skm.mean_absolute_error(y_teste, predicoes_MT )
4 n2 = mean_squared_error(y_teste, predicoes_MT)
5 n3 = skm.max_error(y_teste, predicoes_MT)
6 n4 = skm.median_absolute_error(y_teste, predicoes_MT )
7
8 #metricas Near neighbors
9 in1 = skm.mean_absolute_error(y_teste, neigh_predicoes )
10 in2 = mean_squared_error(y_teste, neigh_predicoes )
11 in3 = skm.max_error(y_teste, neigh_predicoes )
12 in4 = skm.median_absolute_error(y_teste, neigh_predicoes )
13
14 #metricas Decision tree
15 it1 = skm.mean_absolute_error(y_teste, tree_predicoes )
16 it2 = mean_squared_error(y_teste, tree_predicoes )
17 it3 = skm.max_error(y_teste, tree_predicoes )
18 it4 = skm.median_absolute_error(y_teste, tree_predicoes )
19
20 #metricas dummy
21 lu1 = skm.mean_absolute_error(y_teste, dummy_predicoes )
22 lu2 = mean_squared_error(y_teste, dummy_predicoes )
23 lu3 = skm.max_error(y_teste, dummy_predicoes )
24 lu4 = skm.median_absolute_error(y_teste, dummy_predicoes )
25
26 metri = {"modelos":["SVM", "Near Neighbors", "Decision Tree", "Dummy"],"mean_absc
27 modelos_regressao = pd.DataFrame(data = metri)
28 modelos_regressao
29

```


	modelos	mean_absolute_error	mean_squared_error	max_error	median_absolut
0	SVM	66.045554	6604.415306	453.506823	5
1	Near Neighbors	63.782104	6432.605891	701.600000	5
2	Decision Tree	83.781634	11382.141464	792.500000	6