

## ▼ IP-2019.1 data statistics

Save this "2019.1.statistics.ipynb file" to Google Drive and open it with Colab.

```

1 import pylab
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns, numpy as np
5 from google.colab import files
6
7 # to save and download figures of paper
8 global SAVE_FIGS_PAPER
9 SAVE_FIGS_PAPER = True
10 if SAVE_FIGS_PAPER:
11     DPI_resolution = 200
12     width_resolution = 2400
13     height_resolution = 1600
14     !pip install plotly>=4.0.0
15     !wget https://github.com/plotly/orca/releases/download/v1.2.1/orca-1.2.1-x86_64
16     !chmod +x /usr/local/bin/orca
17     !apt-get install xvfb libgtk2.0-0 libgconf-2-4

```

```

--2020-12-25 21:24:00-- https://github.com/plotly/orca/releases/download/v1.2.1
Resolving github.com (github.com)... 140.82.121.3
Connecting to github.com (github.com)|140.82.121.3|:443... connected.
HTTP request sent, awaiting response... 302 Found
Location: https://github-production-release-asset-2e65be.s3.amazonaws.com/990372
--2020-12-25 21:24:01-- https://github-production-release-asset-2e65be.s3.amazo
Resolving github-production-release-asset-2e65be.s3.amazonaws.com (github-produc
Connecting to github-production-release-asset-2e65be.s3.amazonaws.com (github-pr
HTTP request sent, awaiting response... 200 OK
Length: 51607939 (49M) [application/octet-stream]
Saving to: '/usr/local/bin/orca'

```

```

/usr/local/bin/orca 100%[=====>] 49.22M 21.1MB/s in 2.3s

```

```

2020-12-25 21:24:03 (21.1 MB/s) - '/usr/local/bin/orca' saved [51607939/51607939

```

```

Reading package lists... Done

```

Salvo com sucesso

Done

```

newest version (2.24.32-1ubuntu1).
libgconf-2-4 is already the newest version (3.2.6-4ubuntu1).
xvfb is already the newest version (2:1.19.6-1ubuntu4.8).
0 upgraded, 0 newly installed, 0 to remove and 15 not upgraded.

```

## ▼ PE deliveries from ALL students

```

1 file1 = 'http://vision.ufabc.edu.br/MCTest/public/CAE2021/2019.1.PE.csv'
2 df_PE = pd.read_csv(file1, decimal=".", sep=',')
3 df_PE['PE'] = df_PE['PE']*10
4 df_PE.sort_values(by=['Class'], inplace=True)
5 df_PE.sample(5)

```

	Class	ProfT	ProfP	PE	PE1	PE2	PE3	PE4	PE5	PE6	PE7	PE8	PE9	PE10	P
<b>874</b>	32	prof11	prof17	0.0	0	0	0	0	0	0	0	0	0	0	0
<b>130</b>	4	prof01	prof20	0.0	0	0	0	0	0	0	0	0	0	0	0
<b>840</b>	30	prof11	prof02	0.0	0	0	0	0	0	0	0	0	0	0	0
<b>906</b>	33	prof11	prof27	0.0	0	0	0	0	0	0	0	0	0	0	0
<b>1041</b>	37	prof11	prof27	0.0	0	0	0	0	0	0	0	0	0	0	0

```

1 limiar = 0 # WE CONSIDER ONLY MEDIAN LESSONS IN EP > 0
2
3 median = df_PE.groupby(['Class'])['PE'].median()
4 v1 = [str(i) for i in list(median[median>limiar].index.values)]
5 dfVPL_PE = df_PE.query('Class == ' + str(v1))
6 print('VPL-PE:' + str(v1))
7
8 v2 = [str(i) for i in list(median[median<=limiar].index.values)]
9 v3 = v2[-2:] # v3 in IP-BL
10 v2 = v2[:-2] # remove IP-BL
11 dfVPL_NOT_PE = df_PE.query('Class == ' + str(v2))
12 print('VPL-NOT-PE:' + str(v2))
13
14 dfVPL_NOT_BL = df_PE.query('Class == ' + str(v3))
15 print('VPL-BL:' + str(v3))

VPL-PE:['2', '4', '6', '8', '11', '12', '13', '15', '17', '19', '21', '25', '27']
VPL-NOT-PE:['1', '3', '5', '7', '9', '10', '14', '16', '18', '20', '22', '23', '']
VPL-BL:['45', '46']

```

```

1 def draw VPL Used(dfVPL, title):
2     size' ] = (20.0, 8.0)
3
4     ax = sns.boxplot(x="Class", y="PE", data=dfVPL, showmeans=True)
5
6     # Calculate number of obs per group & median to position labels
7     medians = dfVPL.groupby(['Class'])['PE'].median().values
8     nobs = dfVPL.groupby(['Class'])['PE'].count()
9     nobs = [str(x) for x in nobs.tolist()]
10
11     ax.spines["top"].set_visible(False)
12     ax.spines["right"].set_visible(False)

```

Salvo com sucesso

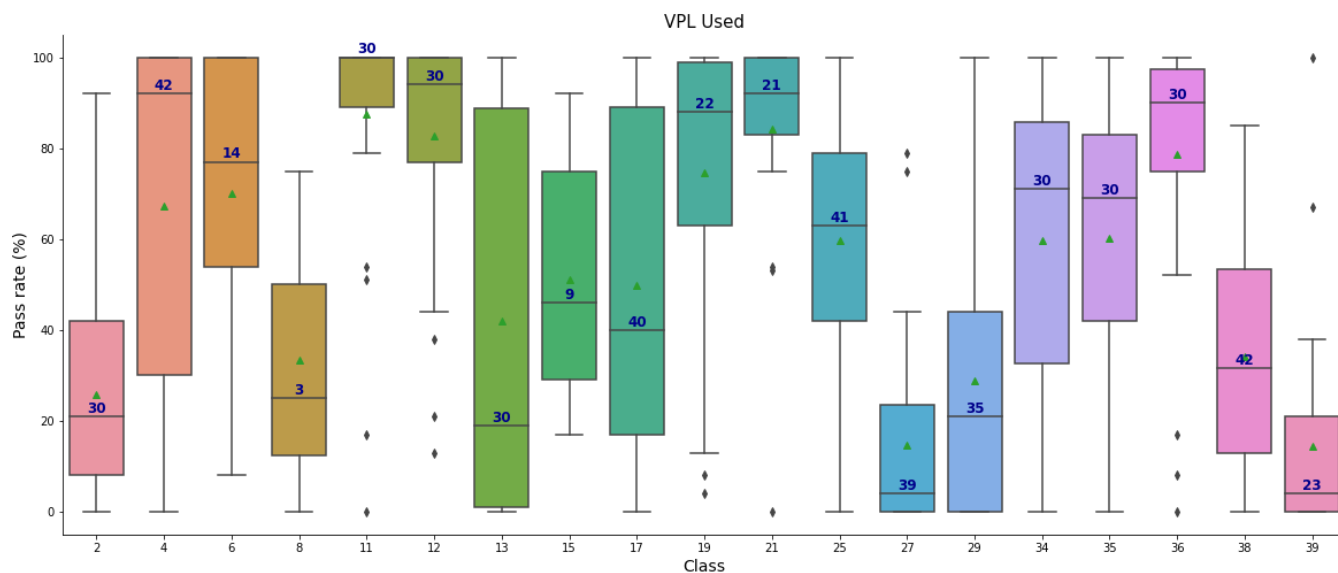


```

12 ax.spines['right'].set_visible(False)
13
14 plt.title(title, size='15')
15
16 pos = range(len(nobs))
17 for tick,label in zip(pos,ax.get_xticklabels()):
18     ax.text(pos[tick], medians[tick] + 0.9, nobs[tick], horizontalalignment='cent
19
20 ax.set_xlabel("Class", fontsize=14)
21 ax.set_ylabel("Pass rate (%)", fontsize=14)
22 ax.plot()

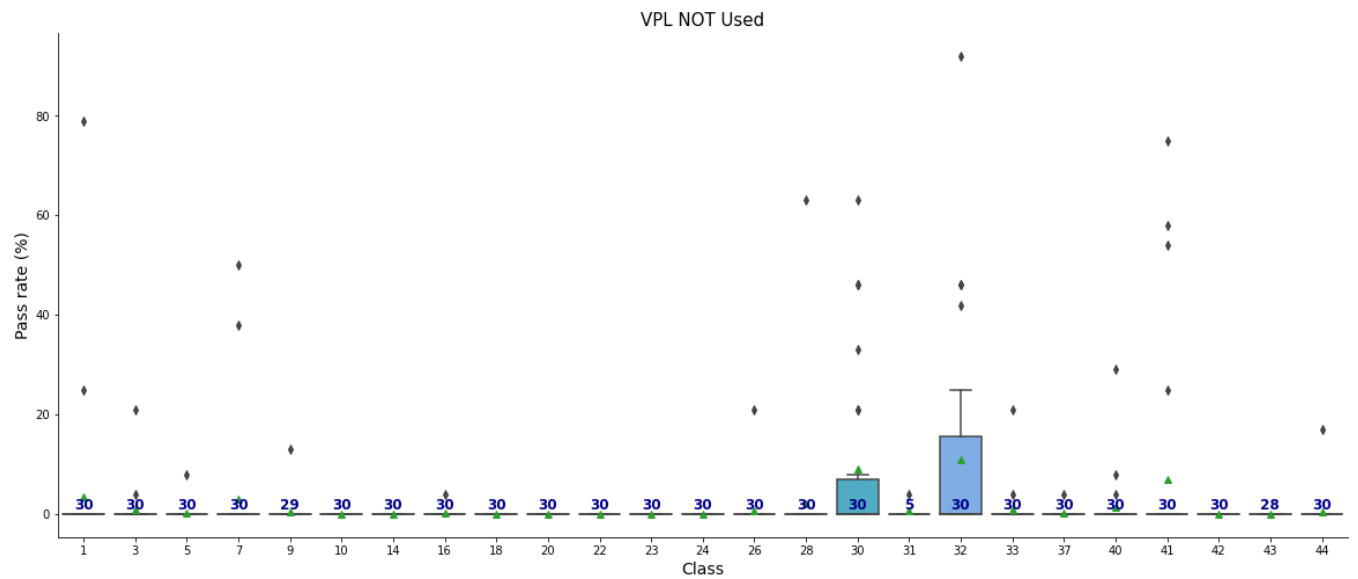
1 title = 'VPL Used'
2 draw_VPL_Used(dfVPL_PE,title)

```

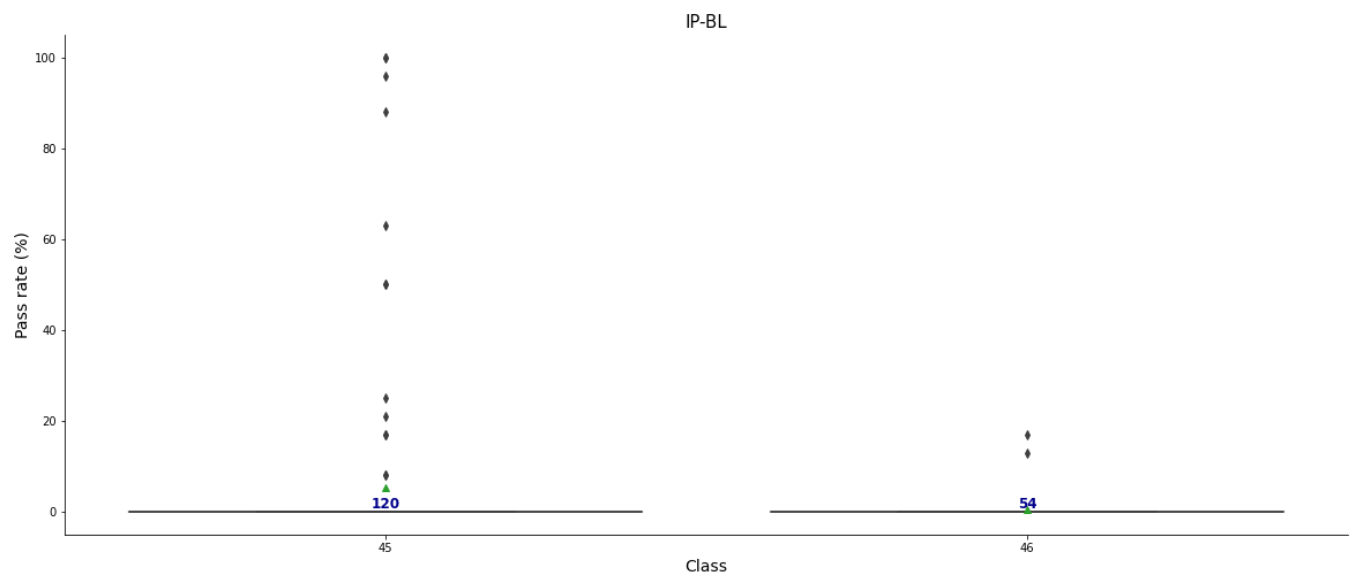


Salvo com sucesso

✕ title)



```
1 title = 'IP-BL'
2 draw_VPL_Used(dfVPL_NOT_BL,title)
```



Salvo com sucesso



In the Boxplot graph, the line in the middle of each rectangle is the median, in red is the number of students who taught some PE.

Cris Sato's explanation of the boxplot:

The box chart presents several useful information:

- The green line inside the box is the median
- The bottom line of the box is the first quartile (Q1), that is, the 25th percentile.
- The top line of the box is the third quartile (Q3), that is, the 75th percentile.
- the lower dash in black is the lowest value in the data above  $Q1 - 1.5IQR$ , where  $IQR = Q3 - Q1$ .
- the upper dash in black is the highest value in the data below  $Q3 + 1.5IQR$ .

The IQR value is called **interquartile range**.

The points below the bottom line or above the top line are called **outliers**.

## ▼ Final grades for ALL 46 classes (academic system)

```

1 file2 = 'http://vision.ufabc.edu.br/MCTest/public/CAE2021/2019.1.grade.csv'
2 df_grade = pd.read_csv(file2, decimal=".")
3 df_grade.sort_values(by=['Adopted','Class'], inplace=True)
4 df_grade['Dropout'] = df_grade['Dropout'].fillna(0)
5 df_grade['New'] = df_grade['New'].fillna(0)
6
7 col_pass_rate = 'Pass rate (%)'
8 df_grade.rename(columns={"Approved": col_pass_rate}, inplace=True)
9
10 df_grade[col_pass_rate] = df_grade[col_pass_rate]*100
11
12 for i in v1:
13     df_grade.loc[df_grade['Class'] == int(i), 'Adopted'] = 'VPL Used'
14
15 for i in v2:
16     df_grade.loc[df_grade['Class'] == int(i), 'Adopted'] = 'VPL NOT Used'
17
18 df_grade.head()

```

Salvo com sucesso



	Class	A	B	C	D	F	O	Total	Dropout	New	Total-Moodle	Total-diff	Pass rate (%)	PE	D/I
0	1	4	5	2	2	6	9	28	2.0	1.0	31	2	46.428571	NaN	I
2	3	3	5	2	3	9	3	25	0.0	0.0	31	6	52.000000	NaN	I

## We include a column with Total-Moodle, including the total number of students enrolled in Moodle

The academic system (PROGRAD) does not offer the exact number of dropouts. We block students and new students in the second week (Adjustment). Thus, abandonment is the official + total difference - Dif2 - Adjustment

```

1 stud_moodle = np.array(df_PE.groupby(['Class'])['PE'].count())
2
3 df_grade.sort_values(by=['Class'], inplace=True)
4 df_grade['Total-Moodle2'] = stud_moodle # what's in moodle - dropout + new
5 df_grade.sort_values(by=['Adopted','Class'], inplace=True)
6
7 df_grade['Total-Diff2'] = df_grade['Total-Moodle2'] - df_grade['Total'] # diff Moodle
8
9 df_grade['Dropout'] += df_grade['Total-Moodle2'] - df_grade['Total'] + df_grade['New']
10 df_grade['Dropout (%)'] = df_grade['Dropout'] / df_grade['Total'] * 100
11
12 df_grade_sort = df_grade.sort_values(by=['Adopted', 'Class'])
13 dfVPL = df_grade_sort.query('Adopted == "VPL Used" and Class != 45 and Class != 46')
14 dfVPL_BL = df_grade_sort.query('Class == 45 or Class == 46')
15 dfVPLNOT = df_grade_sort.query('Adopted == "VPL NOT Used"')
16
17 print("Total-Moodle2 %d " % df_grade['Total-Moodle2'].sum())
18 print("Total-Moodle (new+dropout): %d " % df_grade['Total-Moodle'].sum())
19 print("Total PROGRAD: %d " % df_grade['Total'].sum())
20 print("Dropout: %d " % df_grade['Dropout'].sum())
21 print("New: %d " % df_grade['New'].sum())
22 print("Total-Diff: %d " % df_grade['Total-diff'].sum())
23 print("Total-Diff2: %d " % df_grade['Total-Diff2'].sum())
24 df_grade.sample(5)

```

Salvo com sucesso



```
Total-Moodle2 1437
Total-Moodle (new+dropout): 1466
Total PROGRAD: 1343
Dropout: 148
New: 27
Total-diff: 123
Total-diff2: 94
```

	Class	A	B	C	D	F	O	Total	Dropout	New	Total-Moodle	Total-diff	Pass rate (%)	PE	I
<b>42</b>	43	1	3	2	3	15	3	27	7.0	3.0	31	4	33.333333	NaN	
<b>17</b>	18	2	2	4	3	9	4	24	9.0	1.0	31	6	45.833333	NaN	
<b>23</b>	24	1	2	3	2	18	0	26	6.0	1.0	31	5	30.769231	NaN	
<b>35</b>	36	5	15	6	2	2	0	30	0.0	0.0	30	0	93.333333	7.88	

```
1 def drawGraphs(df_grade, dfVPL, dfVPLNOT, col, file):
2
3     xy = df_grade[['Class',col]]
4     nobs = df_grade[col].values
5
6     # remove valores negativos???
7     if col == 'Dropout (%)':
8         xy[col][xy[col] < 0] = 0
9
10    ax = xy.plot(kind='bar',x='Class',y=col)
11    ax.get_legend().remove()
12
13    x = 24.5
14    x1 = 12
15    x2 = 36.3
16    x3 = 43.5
17
18    plt.axvline(x, color="k", linestyle="--");
19    plt.axvline(x1, color="k", linestyle="--");
20    plt.axvline(x2, color="k", linestyle="--");
21    plt.axvline(x3, color="k", linestyle="--");
22
23    NOTusedMean = dfVPLNOT[col].mean()
24    y1 = NOTusedMean
25    plt.plot([-1, x], [y1, y1], color='red', linewidth=1.5, linestyle="-")
26    #ax.annotate('Avg: %.1f'% y1, xy=(x1, y1+3), size='15', color='red')
27
28    NOTusedSTD = dfVPLNOT[col].std()
29    y1 = NOTusedMean+NOTusedSTD
```

Salvo com sucesso



```

30 plt.plot([-1, x], [y1, y1], color='red', linewidth=1, linestyle="--")
31 #ax.annotate('+Std: %.1f'% y1, xy=(x1, y1+3), size='12', color='red')
32
33 y1 = NOTusedMean-NOTusedSTD
34 plt.plot([-1, x], [y1, y1], color='red', linewidth=1, linestyle="--")
35 #ax.annotate('-Std: %.1f'% y1, xy=(x1, y1+4), size='12', color='red')
36
37 ##### Used
38
39 usedMean = dfVPL[col].mean()
40 y1 = usedMean
41 plt.plot([x, x3], [y1, y1], color='red', linewidth=1.5, linestyle="-")
42 #ax.annotate('Avg: %.1f'% y1, xy=(x2, y1+3), size='15', color='red')
43
44 usedSTD = dfVPL[col].std()
45 y1 = usedMean+usedSTD
46 plt.plot([x, x3], [y1, y1], color='red', linewidth=1, linestyle="--")
47 #ax.annotate('+Std: %.1f'% y1, xy=(x2, y1+3), size='12', color='red')
48
49 y1 = usedMean-usedSTD
50 plt.plot([x, x3], [y1, y1], color='red', linewidth=1, linestyle="--")
51 #ax.annotate('-Std: %.1f'% y1, xy=(x2, y1+3), size='12', color='red')
52
53 y0 = nobs.max() + 0.6
54
55 ax.annotate('VPL NOT Used (Avg: %.1f and Std: %.1f)' % (NOTusedMean, NOTusedSTI
56 ax.annotate('VPL Used (Avg: %.1f and Std: %.1f)' % (usedMean, usedSTD), xy=(29,
57 ax.annotate('IP-BL', xy=(44, y0), size='15')
58 ax.spines["top"].set_visible(False)
59 ax.spines["right"].set_visible(False)
60 #ax.grid(linewidth=.5)
61
62 ax.set_xlabel("Class", fontsize=14)
63 ax.set_ylabel(col, fontsize=14)
64 ax.plot()
65
66 if SAVE_FIGS_PAPER:
67     plt.savefig(file, dpi=DPI_resolution)
68     files.download(file)
69     print('\n\n')

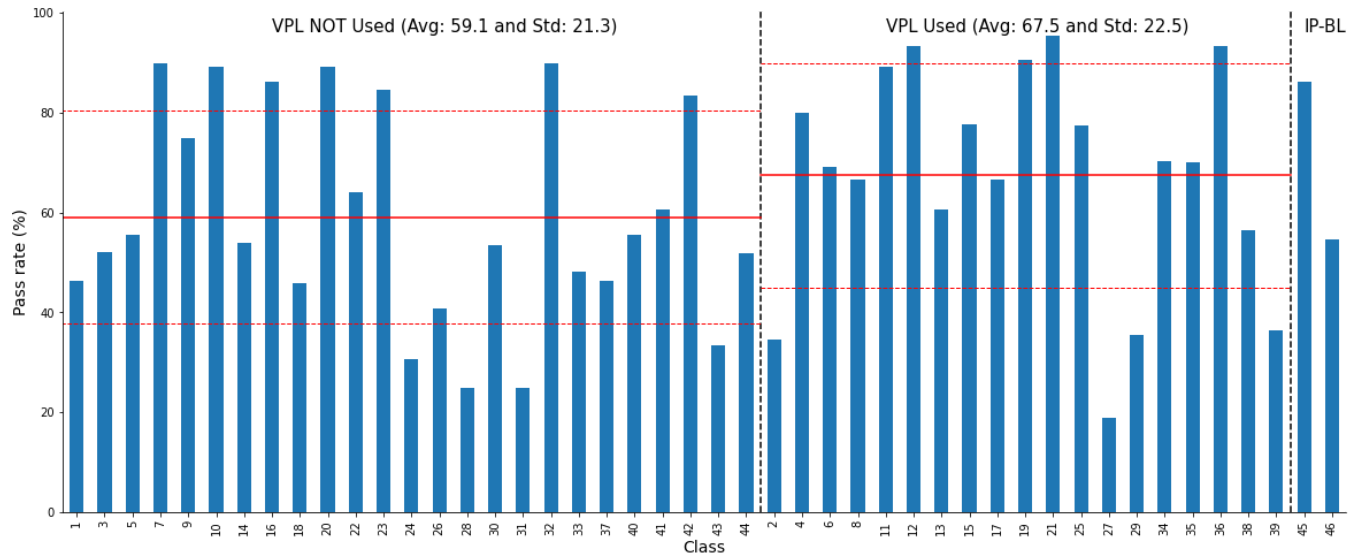
```

Salvo com sucesso



VPLNOT,col\_pass\_rate, 'fig02new.tif')





```
1 drawGraphs(df_grade,dfVPL,dfVPLNOT,'Dropout (%)', 'fig06new.tif')
```

Salvo com sucesso



/usr/local/lib/python3.6/dist-packages/ipykernel\_launcher.py:8: SettingWithCopyWarning: A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: <https://pandas.pydata.org/pandas-docs/stab>

/usr/local/lib/python3.6/dist-packages/pandas/core/series.py:1021: SettingWithCopyWarning: A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: <https://pandas.pydata.org/pandas-docs/stab>

```

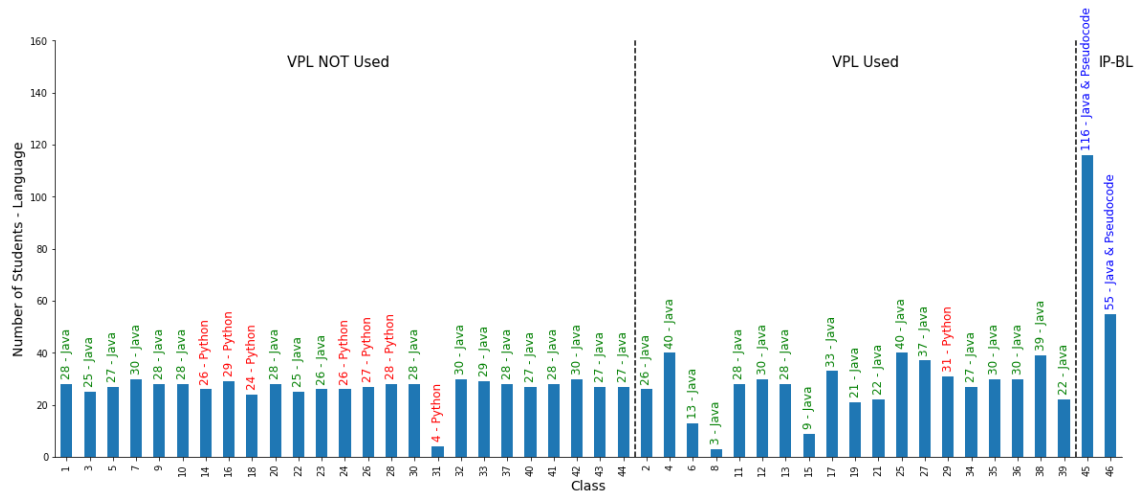
1 xyz = df_grade_sort[['Class', 'Total', 'Language']]
2 ax = xyz.plot(kind='bar', x='Class', y='Total')
3 ax.set_ylim([0, 160])
4 ax.get_legend().remove()
5
6 x = 24.5
7 x1 = 12
8 x2 = 36.3
9 x3 = 43.5
10
11 pos = range(len(xyz['Language']))
12 for tick, label in zip(pos, ax.get_xticklabels()):
13     y0 = xyz['Total'].loc[int(label.get_text())-1]
14     lang0 = xyz['Language'].loc[int(label.get_text())-1]
15     lang = str(y0) + ' - ' + lang0
16     # [ 'normal' | 'bold' | 'heavy' | 'light' | 'ultrabold' | 'ultralight' ]
17     #if int(y)<60:
18     if lang0 == 'Java':
19         ax.text(pos[tick], y0 + 3, lang, horizontalalignment='center', size='12', col
20 elif lang0 == 'Java & Pseudocode':
21         ax.text(pos[tick], y0 + 3, lang, horizontalalignment='center', size='12', col
22 else:
23         ax.text(pos[tick], y0 + 3, lang, horizontalalignment='center', size='12', col
24
25
26 plt.axvline(x, color="k", linestyle="--");
27 plt.axvline(x3, color="k", linestyle="--");
28
29 ax.annotate('VPL NOT Used', xy=(9.5, 150), size='15')
30 ax.annotate('VPL Used', xy=(33, 150), size='15')
31 ax.annotate('IP-BL', xy=(44.5, 150), size='15')
32
33
34
35
36 ax.set_xlabel("Class", fontsize=14)
37 ax.set_ylabel("Number of Students - Language", fontsize=14)
38
39 if SAVE_FIGS_PAPER:
40     plt.savefig('fig05new.tif', dpi=DPI_resolution)
41     files.download('fig05new.tif')

```

Salvo com sucesso



le(False)



```
1 mean = dfVPL.groupby(['Class'])['PE'].mean()
2 nobs = dfVPL.groupby(['Class'])['PE'].count()
```

```
1 col,lin = df_grade.shape
2 x = []
3 for i in range(0,col):
4     v = df_grade[['Class',col_pass_rate]].loc[i]
5     try:
6         x.append([i+1,v[1],round(mean[v[0]],1),nobs[v[0]])]
7     except:
8         pass
```

Salvo com sucesso

```
array([[ 2.      , 34.615385,  2.6      ,  1.      ],
       [ 4.      ,  80.      ,  6.8      ,  1.      ],
       [ 6.      , 69.230769,  7.      ,  1.      ],
       [ 8.      , 66.666667,  3.3      ,  1.      ],
       [11.     , 89.285714,  8.8      ,  1.      ],
       [12.     , 93.333333,  8.3      ,  1.      ],
       [13.     , 60.714286,  4.2      ,  1.      ],
       [15.     , 77.777778,  5.2      ,  1.      ],
       [17.     , 66.666667,  5.      ,  1.      ]],
```

```
[19.      , 90.47619 , 7.1      , 1.      ],
[21.      , 95.454545, 8.4      , 1.      ],
[25.      , 77.5      , 6.      , 1.      ],
[27.      , 18.918919, 1.5      , 1.      ],
[29.      , 35.483871, 2.9      , 1.      ],
[34.      , 70.37037 , 6.      , 1.      ],
[35.      , 70.      , 6.      , 1.      ],
[36.      , 93.333333, 7.9      , 1.      ],
[38.      , 56.410256, 3.4      , 1.      ],
[39.      , 36.363636, 1.4      , 1.      ]])
```

```
1 data = {'Class': x[:,0].astype(int),
2         'Pass rate': x[:,1],
3         'PE': x[:,2],
4         'N' : x[:,3].astype(int)}
5
6 grade_vpl = pd.DataFrame(data)
```

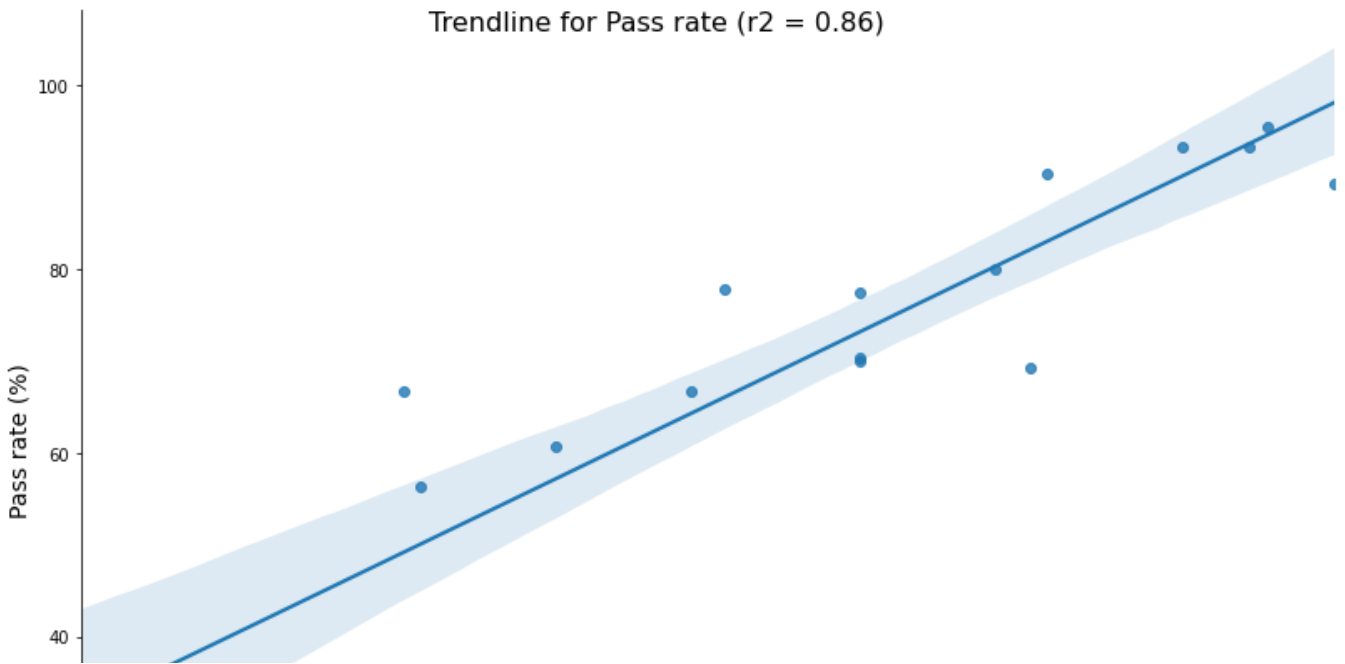
```
1 correl = grade_vpl['Pass rate'].corr(grade_vpl['PE'])
2 r2 = correl**2
3 print(correl, r2)
```

```
0.9293857826418399 0.8637579329767853
```

```
1 import seaborn as sns
2 x=sns.lmplot(x="PE", y="Pass rate", data=grade_vpl)
3 x.fig.set_size_inches(12,8)
4 x.fig.suptitle('Trendline for Pass rate (r2 = %.2f)'%r2, fontsize=16, y=0.97)
5 x.set_xlabel("Average Grade of VPL Programming Exercises (PE) (%)", fontsize=14)
6 x.set_ylabel("Pass rate (%)", fontsize=14)
7 #x.set_xlim(0,100)
8 #x.set_xticks(range(1,32))
9
10 if SAVE_FIGS_PAPER:
11     plt.savefig('fig04new.tif', dpi=DPI_resolution)
12     files.download('fig04new.tif')
```

Salvo com sucesso



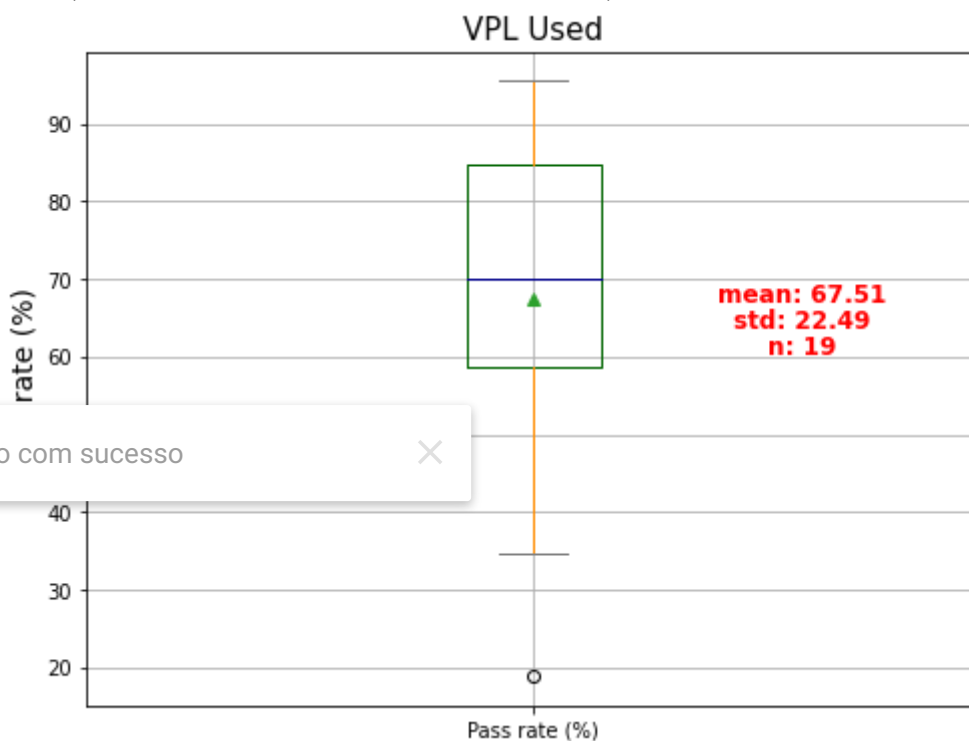


```

1 color = {'boxes': 'DarkGreen', 'whiskers': 'DarkOrange', 'medians': 'DarkBlue', 'color': 'DarkBlue'}
2 ax = dfVPL.plot.box(y=col_pass_rate, figsize=(8, 6), showmeans=True, color=color)
3 ax.grid()
4 plt.title('VPL Used', size='15')
5
6 ax.set_ylabel(col_pass_rate, fontsize=14)
7
8 ax.text(1.3, dfVPL[col_pass_rate].mean()-0.5, 'mean: '+str(round(dfVPL[col_pass_rate].mean(), 2)))
9 ax.text(1.3, dfVPL[col_pass_rate].mean()-3.9, 'std: '+str(round(dfVPL[col_pass_rate].std(), 2)))
10 ax.text(1.3, dfVPL[col_pass_rate].mean()-7.2, 'n: '+str(dfVPL[col_pass_rate].count()))

```

Text(1.3, 60.30535363157894, 'n: 19')

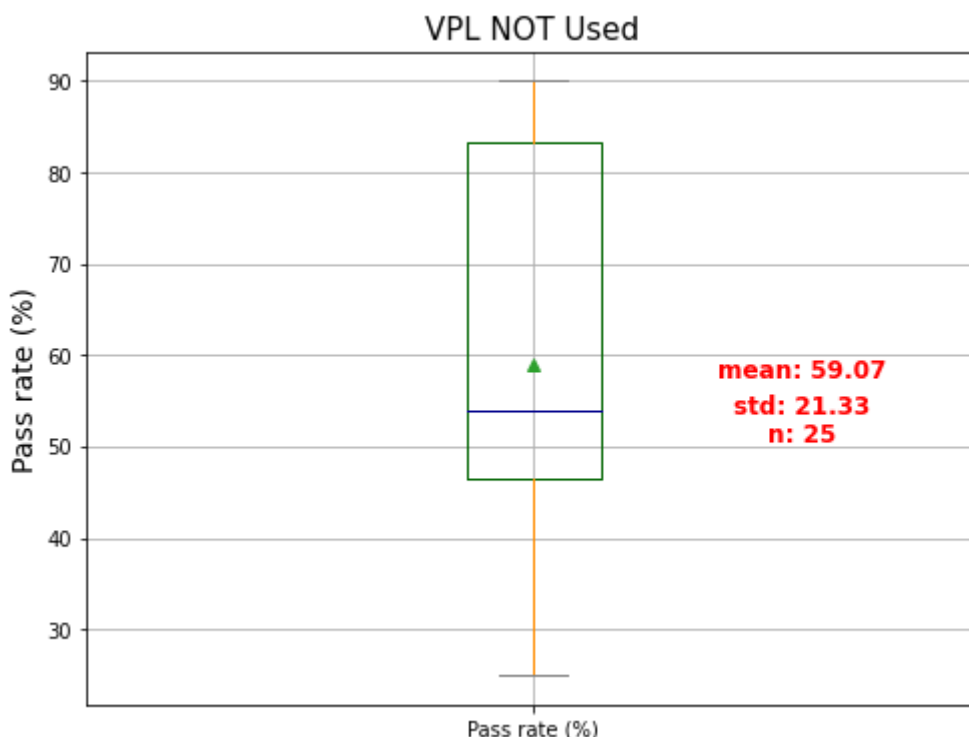


```

1 color = {'boxes': 'DarkGreen', 'whiskers': 'DarkOrange', 'medians': 'DarkBlue', 'c
2 ax = dfVPLNOT.plot.box(y=col_pass_rate, figsize=(8, 6), showmeans=True, color=col
3 ax.grid()
4 plt.title('VPL NOT Used', size='15')
5 ax.set_ylabel(col_pass_rate, fontsize=14)
6 ax.text(1.3, dfVPLNOT[col_pass_rate].mean()-1.5, 'mean: '+str(round(dfVPLNOT[col_
7 ax.text(1.3, dfVPLNOT[col_pass_rate].mean()-5.5, 'std: '+str(round(dfVPLNOT[col_r
8 ax.text(1.3, dfVPLNOT[col_pass_rate].mean()-8.5, 'n: '+str(dfVPLNOT[col_pass_rate
9

```

Text(1.3, 50.56526072, 'n: 25')



```

1 import plotly.graph_objects as go
2 from plotly.subplots import make_subplots
3
4 fig1 = go.Box(
5     y=dfVPLNOT[col_pass_rate].values,
6     name='NOT Used',
7     boxmean='sd'
8 )
9 fig2 = go.Box(
10    y=dfVPLNOT[col_pass_rate].values,
11    name='Used',
12    boxmean='sd'
13 )
14
15 #boxfig= go.Figure(data=[fig1,fig2])
16 #fig = make_subplots(rows=1, cols=2)
17 #fig.add_trace(boxfig.data[0], row=1, col=1)
18 #fig.add_trace(boxfig.data[1], row=1, col=2)
19

```

Salvo com sucesso

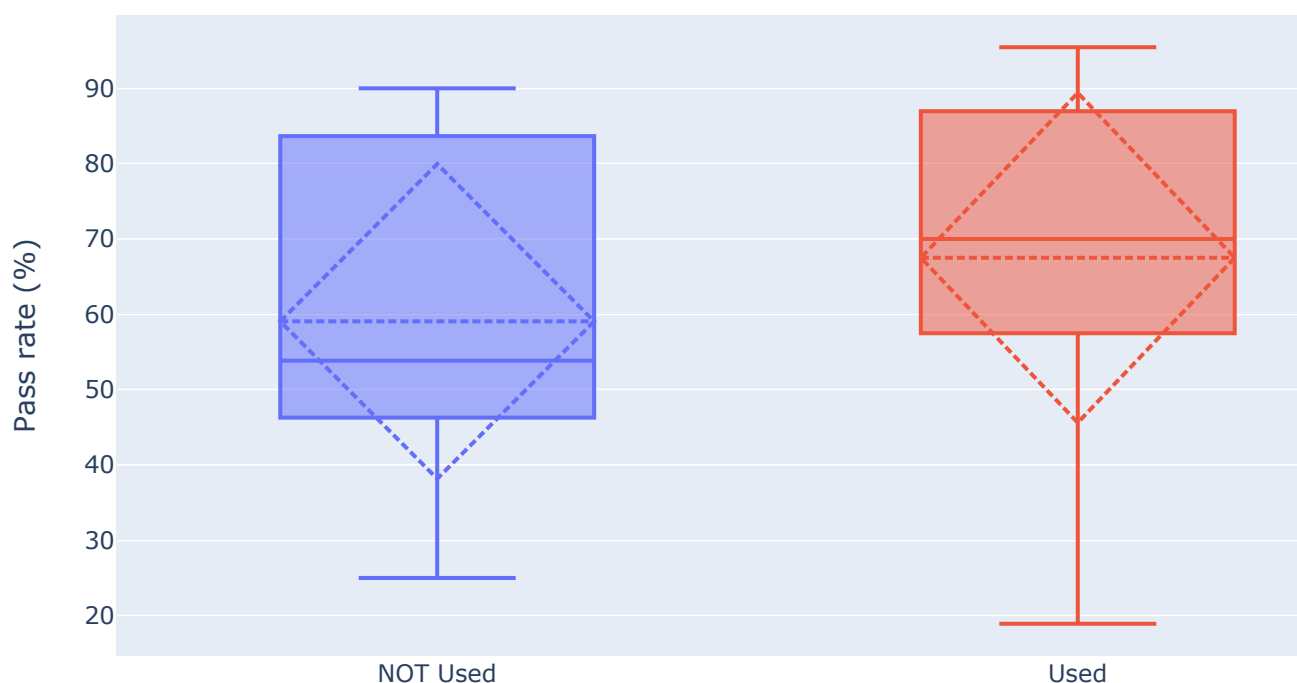
values,

```

20 layout = go.Layout(
21     width=800,height=500,showlegend=False,
22     title={'text': "BoxPlot",'y':0.9,'x':0.5,'xanchor': 'center','yanchor': 'top'}
23     yaxis=dict(
24         title=col_pass_rate
25     ) )
26 fig = go.Figure(layout=layout)
27 fig.add_trace(fig1)
28 fig.add_trace(fig2)
29 fig.show()

```

BoxPlot



```

1 import plotly.graph_objects as go
2 fig1 = go.Violin(meanline_visible=True,box_visible = True,
3                 y=dfVPL[col_pass_rate].values,

```

Salvo com sucesso

```

6 fig2 = go.Violin(meanline_visible=True,box_visible = True,
7                 y=dfVPL[col_pass_rate].values,
8                 name='Used'
9             )
10 #style layout
11 layout = go.Layout(
12     width=800,height=500,showlegend=False,
13     title={'text': "BoxPlot & Violin Plots",'y':0.9,'x':0.5,'xanchor': 'center','

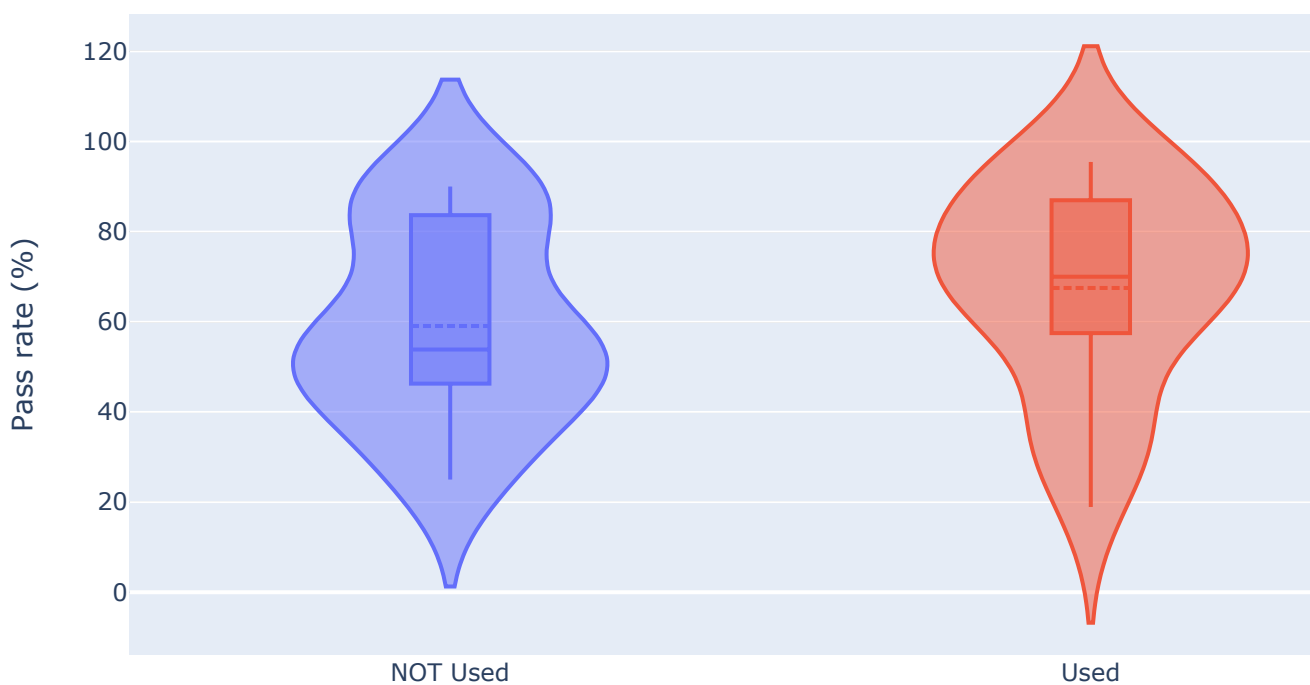
```

```

14     yaxis=dict(
15         title=col_pass_rate
16     ) )
17 fig = go.Figure(layout=layout)
18 fig.add_trace(fig1)
19 fig.add_trace(fig2)
20 fig.show()
21
22 # kde(kernel density estimation).

```

## BoxPlot & Violin Plots



```

1 import plotly.graph_objects as go
2 import matplotlib.pyplot as plt
3
4 fig1 = go.Violin(meanline_visible=True,box_visible=True,line_color='black',opacit
5     y=dfVPL[NOT[col_pass_rate]],
6     fillcolor='blueviolet',
7     name='NOT Used'
8 )
9
10 fig2 = go.Violin(meanline_visible=True,box_visible=True,line_color='black',opacit
11     y=dfVPL[col_pass_rate],
12     fillcolor='orangered',
13     name='Used'
14 )

```

Salvo com sucesso

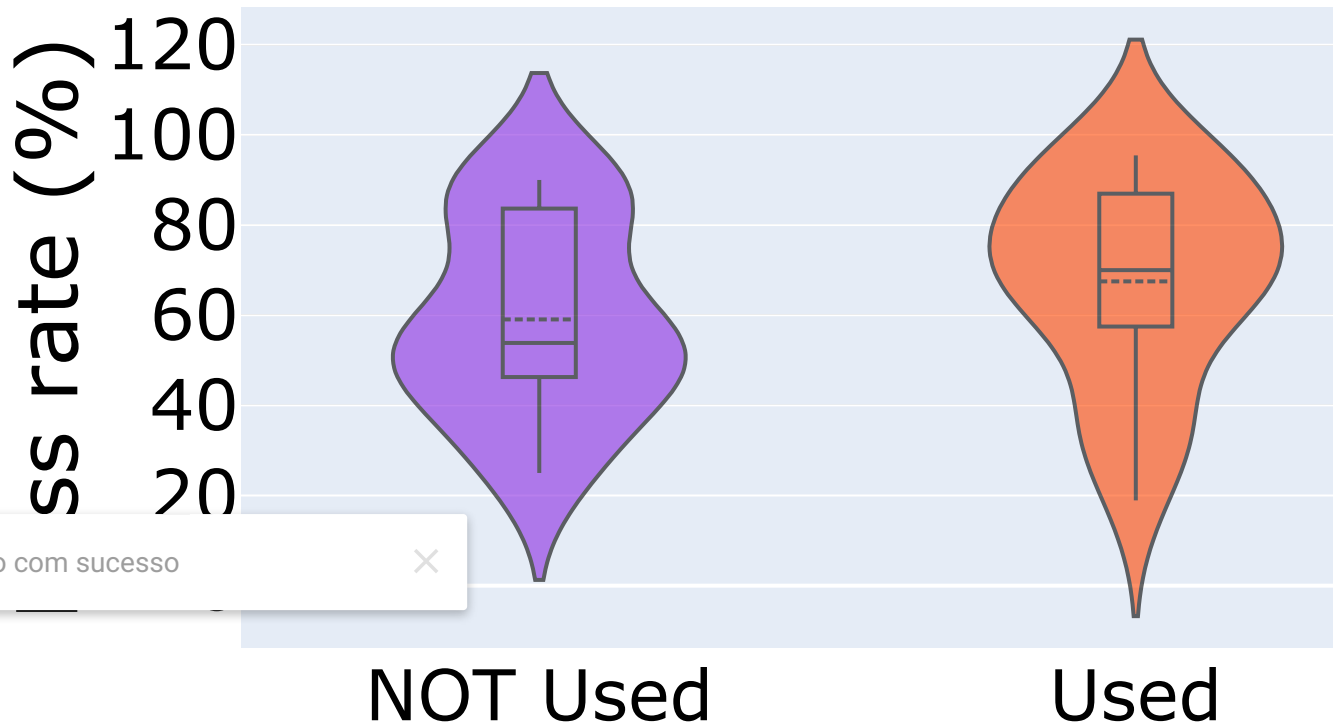


visible=True,box\_visible=True,line\_color='black',opacit



```
15 ,
16 #style layout
17 layout = go.Layout(title_font_size=50,
18     width=800,height=500,showlegend=False,
19     title={'text': "BoxPlot & Violin Plots",'y':0.9,'x':0.5,'xanchor': 'center','
20     font=dict(size=35, color = 'black'),
21     yaxis=dict(
22         title=col_pass_rate,
23     ) )
24
25 fig = go.Figure(layout=layout)
26 fig.add_trace(fig1)
27 fig.add_trace(fig2)
28 fig.show()
29
30 if SAVE_FIGS_PAPER:
31     from PIL import Image
32     fig.write_image('fig03.png', width=width_resolution, height=height_resolution)
33     img = Image.open('fig03.png')
34     img.save('fig03.tif')
35     files.download('fig03.tif')
```

# BoxPlot & Violin Plots



Salvo com sucesso

## Refs

[boxplot] Tukey, John W. Exploratory data analysis. Mosteller F, editor. United States Addison-Wesley. Publishing Company; (1977): 688 p.

[ViolinPlots] Hintze, Jerry L., and Ray D. Nelson. Violin plots: a box plot-density trace synergism. The American Statistician. 52.2 (1998): 181-184.

## Cores:

aliceblue, antiquewhite, aqua, aquamarine, azure,  
 beige, bisque, black, blanchedalmond, blue,  
 blueviolet, brown, burlywood, cadetblue,  
 chartreuse, chocolate, coral, cornflowerblue,  
 cornsilk, crimson, cyan, darkblue, darkcyan,  
 darkgoldenrod, darkgray, darkgrey, darkgreen,  
 darkkhaki, darkmagenta, darkolivegreen, darkorange,  
 darkorchid, darkred, darksalmon, darkseagreen,  
 darkslateblue, darkslategray, darkslategrey,  
 darkturquoise, darkviolet, deeppink, deepskyblue,  
 dimgray, dimgrey, dodgerblue, firebrick,  
 floralwhite, forestgreen, fuchsia, gainsboro,  
 ghostwhite, gold, goldenrod, gray, grey, green,  
 greenyellow, honeydew, hotpink, indianred, indigo,  
 ivory, khaki, lavender, lavenderblush, lawngreen,  
 lemonchiffon, lightblue, lightcoral, lightcyan,  
 lightgoldenrodyellow, lightgray, lightgrey,  
 lightgreen, lightpink, lightsalmon, lightseagreen,  
 lightskyblue, lightslategray, lightslategrey,  
 lightsteelblue, lightyellow, lime, limegreen,  
 linen, magenta, maroon, mediumaquamarine,  
 mediumblue, mediumorchid, mediumpurple,  
 mediumseagreen, mediumslateblue, mediumspringgreen,  
 mediumturquoise, mediumvioletred, midnightblue,  
 mintcream, mistyrose, moccasin, navajowhite, navy,  
 orange, orangered,  
 green, paleturquoise,  
 palevioletred, papayawhip, peachpuff, peru, pink,  
 plum, powderblue, purple, red, rosybrown,  
 royalblue, rebeccapurple, saddlebrown, salmon,  
 sandybrown, seagreen, seashell, sienna, silver,  
 skyblue, slateblue, slategray, slategrey, snow,  
 springgreen, steelblue, tan, teal, thistle, tomato,

Salvo com sucesso



```
turquoise, violet, wheat, white, whitesmoke,
yellow, yellowgreen
```

## ▼ Students' T Test

```
1 countVPL = dfVPL[['A','B','C','D','F','O']].sum().sum()
2 countVPLNOT = dfVPLNOT[['A','B','C','D','F','O']].sum().sum()
3 countVPL_BL = dfVPL_BL[['A','B','C','D','F','O']].sum().sum()
4 meanVPL = dfVPL[col_pass_rate].mean()
5 meanVPLNOT = dfVPLNOT[col_pass_rate].mean()
6 meanVPL_BL = dfVPL_BL[col_pass_rate].mean()
7 stdVPL = dfVPL[col_pass_rate].std()
8 stdVPLNOT = dfVPLNOT[col_pass_rate].std()
9 stdVPL_BL = dfVPL_BL[col_pass_rate].std()
10
```

```
1 print(countVPL_BL, meanVPL_BL)
2 print(countVPL, meanVPL)
3 print(countVPLNOT, meanVPLNOT)
4
```

```
171 70.376176
509 67.50535363157894
663 59.06526072
```

```
1 # calculate means
2 mean1, mean2 = meanVPL, meanVPLNOT
```

```
1 # calculate sample standard deviations
2 std1, std2 = stdVPL, stdVPLNOT
```

```
1 # calculate standard errors
2 import math
3 n1, n2 = countVPL, countVPLNOT
4 se1, se2 = std1/math.sqrt(n1), std2/math.sqrt(n2)
```

```
Salvo com sucesso × Difference between the samples
2 t_stat = (mean1 - mean2) / (se1 + se2**2.0)
```

```
1 # calculate the t statistic
2 t_stat = (mean1 - mean2) / se
3 t_stat
```

```
6.512378158662304
```

```

1 # degrees of freedom
2 df = n1 + n2 - 2
3 df

```

```
1170
```

```

1 from scipy.stats import t
2 # calculate the critical value
3 alpha = 0.05
4 cv = t.ppf(1.0 - alpha, df)
5 cv

```

```
1.6461570318357832
```

```

1 # calculate the p-value
2 p = (1.0 - t.cdf(abs(t_stat), df)) * 2.0
3 '%.15f' % p

```

```
'0.000000000109572'
```

```

1 # confidence level
2 '%.10f' % ((1-p)*100)

```

```
'99.9999999890'
```

## ▼ Similarity in PE

```

1 file3 = 'http://vision.ufabc.edu.br/MCTest/public/CAE2021/2019.1.similarity.csv'
2 df_similarity = pd.read_csv(file3)
3 df_similarity.sort_values(by=['PE'], inplace=True)
4 df_similarity['Similarity (%)'] = df_similarity.Similarity / df_similarity.Subscr
5 df_similarity.head(5)

```

↳

	PE	Similarity	Subscriptions	Corrections	Type	Similarity (%)
0	1	2	596	582	Sequence	0.335570
			547	506	Sequence	26.142596
2	3	158	499	468	Sequence	31.663327
3	4	145	483	454	Sequence	30.020704
4	5	136	452	422	Sequence	30.088496

```

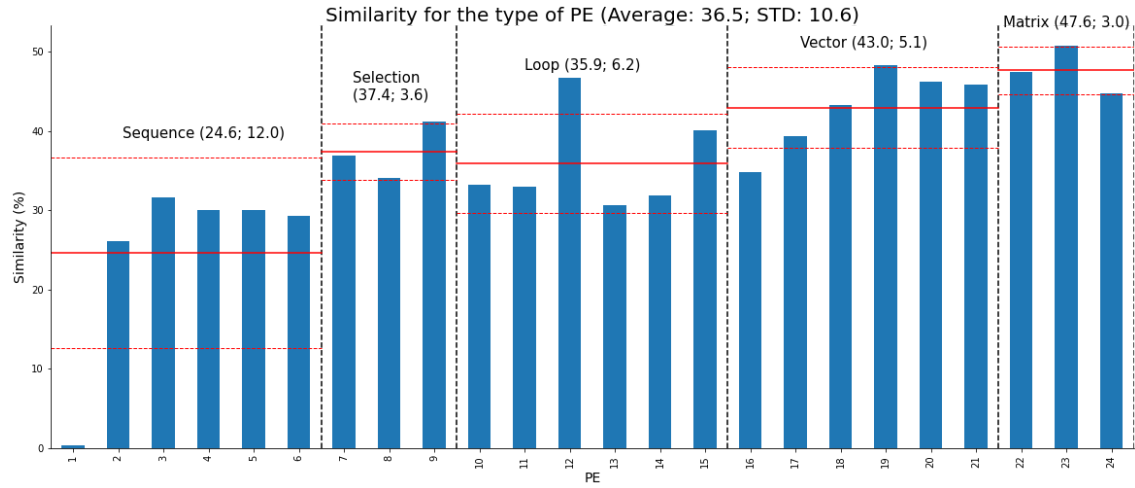
1 column = 'Similarity (%)'
2 xy = df_similarity[['PE', column]]
3 ax = xv.plot(kind='bar', x='PE', v=column)

```

```
3 ax.get_legend().remove()
4 ax.spines["top"].set_visible(False)
5 ax.spines["right"].set_visible(False)
6
7
8 types = ['Sequence', 'Selection', 'Loop', 'Vector', 'Matrix']
9
10 col_sum = 0
11 for t in types:
12     df_aux = df_similarity.query('Type == "' + t + '"')
13     mean = df_aux[column].mean()
14     std = df_aux[column].std()
15     dx, _ = df_aux.shape
16     col_sum += dx
17     x1 = col_sum - .5
18     x0 = x1 - dx
19     plt.axvline(x1, color="k", linestyle="--");
20     plt.plot([x0, x1], [mean, mean], color='red', linewidth=1.5, linestyle="-")
21     plt.plot([x0, x1], [mean-std, mean-std], color='red', linewidth=1, linestyle="-")
22     plt.plot([x0, x1], [mean+std, mean+std], color='red', linewidth=1, linestyle="-")
23     if t == 'Selection':
24         ax.annotate('%s \n(%.1f; %.1f)' % (t, mean, std), xy=((x0+x1)/2-0.8, mean+std))
25     elif t == 'Loop':
26         ax.annotate('%s (%.1f; %.1f)' % (t, mean, std), xy=((x0+x1)/2-1.5, mean+std +
27     else:
28         ax.annotate('%s (%.1f; %.1f)' % (t, mean, std), xy=((x0+x1)/2-1.4, mean+std +
29
30 mean = df_similarity['Similarity (%)'].mean()
31 std = df_similarity['Similarity (%)'].std()
32 #print("mean:",str(mean),"std:",str(std))
33 plt.title('Similarity for the type of PE (Average: %.1f; STD: %.1f)'%(mean,std),
34 ax.set_xlabel("PE", fontsize=14)
35 ax.set_ylabel(column, fontsize=14)
36 ax.plot()
37
38 if SAVE_FIGS_PAPER:
39     plt.savefig('fig07new.tif', dpi=DPI_resolution)
40     files.download('fig07new.tif')
41 print('\n\n')
```

Salvo com sucesso





```
1 df_similarity.sample(3)
```

PE	Similarity	Subscriptions	Corrections	Type	Similarity (%)	
7	8	137	402	387	Selection	34.079602
6	7	153	414	397	Selection	36.956522
0	1	2	596	582	Sequence	0.335570

1

Salvo com sucesso ✕

Salvo com sucesso

