```
1 pip install pandas-profiling
```

```
  Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheel
  Collecting pandas-profiling
    Downloading pandas-profiling-3.6.6-py2.py3-none-any.whl (324 kB)
```

```
 1 import pandas as pd
 2 import io
 3 import matplotlib
 4 import scipy
 5 from scipy import stats as stats
 6 import seaborn as sns
 7 import matplotlib.pyplot as plt
 8 import math
 9 import numpy as np
10 from pandas_profiling import ProfileReport
11
```

```
    <ipython-input-2-9cb39e080237>:10: DeprecationWarning: `import pandas_profiling` is going to be deprecated by April 1st. Please use
      from pandas_profiling import ProfileReport
```

```
    Downloading visions-0.7.5-py3-none-any.whl (102 kB)
```

```
 1 test_data = pd.read_csv('test_processed.csv')
 2
 3 ProfileReport(test_data, minimal=True)
```

| Summarize dataset: | 50/50 [00:01<00:00, 48.79it/s, |
|---|---|
| 100% | Completed] |
| Generate report structure: | 1/1 [01:11<00:00, |
| 100% | 71.17s/it] |
| Render HTML: | 1/1 [00:01<00:00, |
| 100% | 1.46s/it] |

# Overview

### Dataset statistics

| | |
|---|---|
| **Number of variables** | 44 |
| **Number of observations** | 10000 |
| **Missing cells** | 0 |
| **Missing cells (%)** | 0.0% |
| **Total size in memory** | 3.4 MiB |
| **Average record size in memory** | 352.0 B |

### Variable types

| | |
|---|---|
| **Numeric** | 44 |

### Alerts

| | |
|---|---|
| `FiO2` is highly skewed ($\gamma 1$ = -43.1230073) | **Skewed** |
| `Alkalinephos` is highly skewed ($\gamma 1$ = 23.28169282) | **Skewed** |

```
 1 train_data = pd.read_csv('train_processed.csv')
 2
 3 ProfileReport(train_data, minimal=True)
```

| | |
|---|---|
| Summarize dataset: | 50/50 [00:00<00:00, 54.63it/s, |
| 100% | Completed] |
| Generate report structure: | 1/1 [00:33<00:00, |
| 100% | 33.68s/it] |
| Render HTML: | 1/1 [00:01<00:00, |
| 100% | 1.27s/it] |

# Overview

### Dataset statistics

| | |
|---|---|
| **Number of variables** | 44 |
| **Number of observations** | 20000 |
| **Missing cells** | 0 |
| **Missing cells (%)** | 0.0% |
| **Total size in memory** | 6.7 MiB |
| **Average record size in memory** | 352.0 B |

### Variable types

| | |
|---|---|
| **Numeric** | 44 |

### Alerts

| | |
|---|---|
| `FiO2` is highly skewed ($\gamma1 = -113.2576163$) | **Skewed** |
| `Bilirubin_direct` is highly skewed ($\gamma1 = 25.4089327$) | **Skewed** |

```
1 sns.catplot(x="SepsisLabel", kind="count", palette="ch:.50", data=test_data).set(title='Sepsis Label')
```

```
<seaborn.axisgrid.FacetGrid at 0x7fd3169ede40>
```
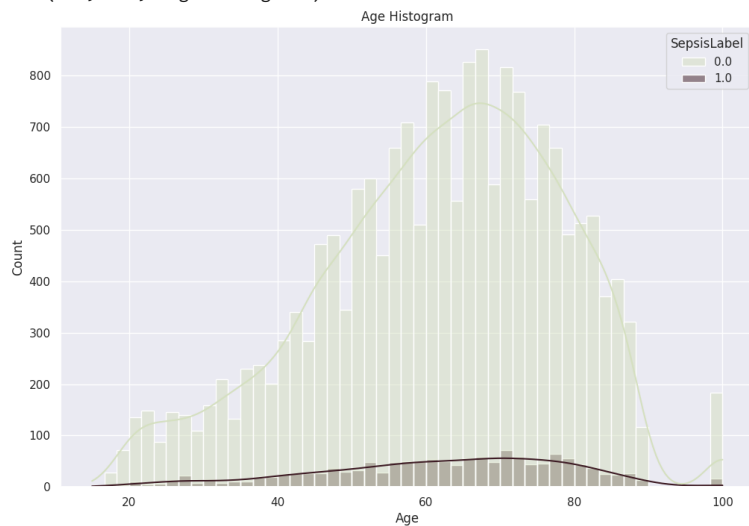


```
1 sns.set(rc={'figure.figsize':(12,8)})
2 sns.histplot(data=test_data, x="Age" ,kde=True, hue="SepsisLabel",  palette="ch:.60").set_title('Age Histogram')
3
```

Text(0.5, 1.0, 'Age Histogram')

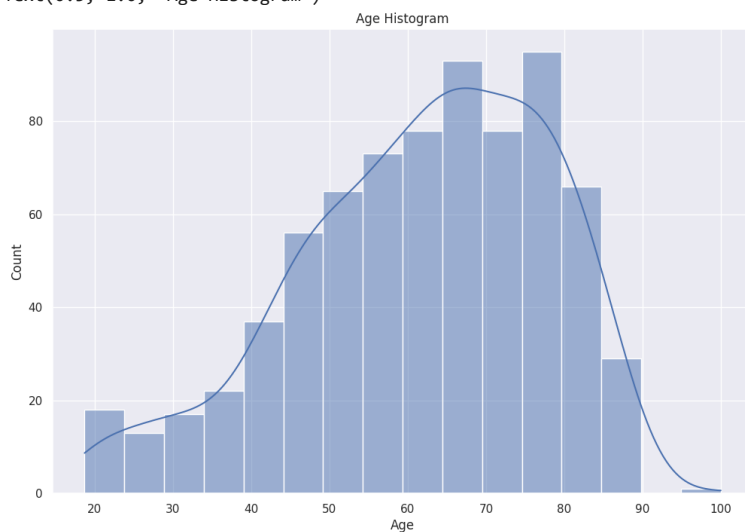

```
1 sns.histplot(data=train_data, x="Age" ,kde=True, hue="SepsisLabel",  palette="ch:.60").set_title('Age Histogram')
2
```

Text(0.5, 1.0, 'Age Histogram')

```
1 sns.histplot(data=test_data[test_data['SepsisLabel']==1], x="Age",kde=True,  palette="hls2").set_title('Age Histogram
2
```
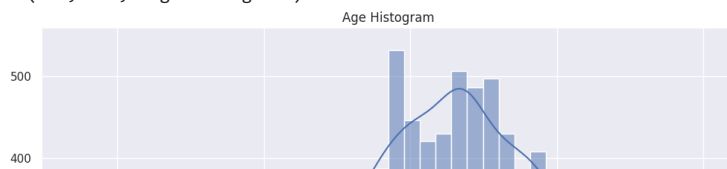
```
<ipython-input-8-3a3226cb0b72>:1: UserWarning: Ignoring `palette` because
  sns.histplot(data=test_data[test_data['SepsisLabel']==1], x="Age",kde=Tr
Text(0.5, 1.0, 'Age Histogram')
```



```
1
```

```
1 sns.histplot(data=test_data[test_data['SepsisLabel']==0], x="Age",kde=True,  palette="hls2").set_title('Age Histogram
2
```

```
<ipython-input-8-3a3226cb0b72>:1: UserWarning: Ignoring `palette` because
  sns.histplot(data=test_data[test_data['SepsisLabel']==1], x="Age",kde=Tr
Text(0.5, 1.0, 'Age Histogram')
```
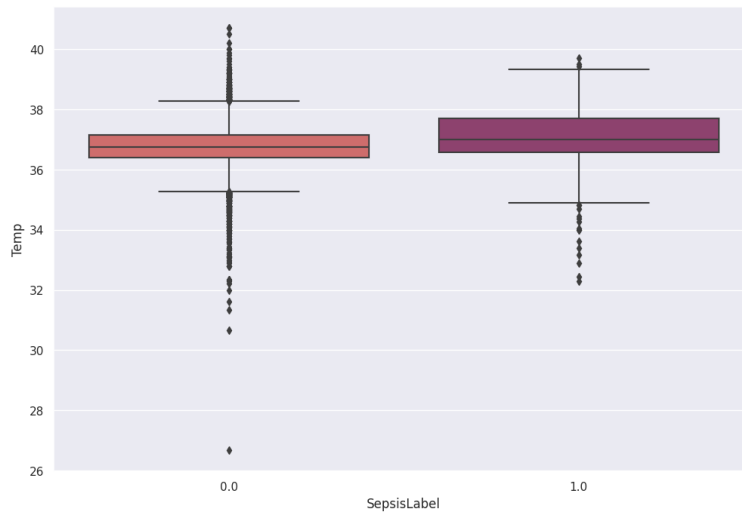
```
<ipython-input-9-880143515de9>:1: UserWarning: Ignoring `palette` because
  sns.histplot(data=test_data[test_data['SepsisLabel']==0], x="Age",kde=Tr
Text(0.5, 1.0, 'Age Histogram')
```
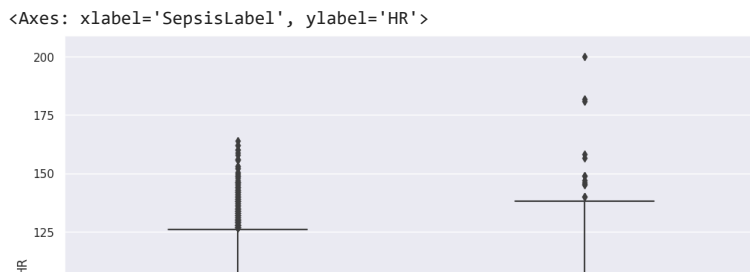


```
1 sns.boxplot(x =train_data["SepsisLabel"], y=train_data["Temp"],  palette ="flare" )
2
```

```
<Axes: xlabel='SepsisLabel', ylabel='Temp'>
```



```
1 sns.boxplot(x =train_data["SepsisLabel"], y=train_data["HR"],  palette ="flare" )
2
```

`<Axes: xlabel='SepsisLabel', ylabel='HR'>`



## מבחנים סטטיסטיים

נבחן האם סימפטומים של אלח דם הם חום גבוה, דופק גבוה,ירידה בטסיות הדם יתכן גם לחץ דם נמוך

```
1 # T test   תחת הנחת התפלגות נורמלית, נראה גם שהשונות דומה ותחת הנחה נוספת שהתפלגות ההפרש היא נורמלית נבצע
```

```python
1
2 sepsis_temp = train_data[train_data['SepsisLabel'] == 1]
3 sepsis_temp = sepsis_temp.Temp
4 sepsis_temp = np.array(sepsis_temp)
5
6
7 not_sepsis_temp = train_data[train_data['SepsisLabel'] == 0]
8 not_sepsis_temp = not_sepsis_temp.Temp
9 not_sepsis_temp = np.array(not_sepsis_temp)
10
11 mu_sepsis_temp = np.mean(sepsis_temp)
12 mu_non_sepsis_temp = np.mean(not_sepsis_temp)
13
14 n_sepsis = len(sepsis_temp)
15 n_non_sepsis = len(not_sepsis_temp)
16
17 sigma_sepsis_temp = math.sqrt(np.sum((sepsis_temp - mu_sepsis_temp)**2) / n_sepsis)
18 sigma_not_sepsis_temp = math.sqrt(np.sum((not_sepsis_temp - mu_non_sepsis_temp)**2) / n_non_sepsis)
19
20 d = mu_sepsis_temp - mu_non_sepsis_temp
21
22 print('mean of temp of patients with sepsis:', mu_sepsis_temp)
23 print('varience of temp of patients with sepsis:', sigma_sepsis_temp)
24 print('mean of temp of patients without sepsis:', mu_non_sepsis_temp)
25 print('varience of temp of patients without sepsis:', sigma_not_sepsis_temp)
26 print('the difference of means is: ', d)
```

```
mean of temp of patients with sepsis: 37.07611741190975
varience of temp of patients with sepsis: 0.873972987759184
mean of temp of patients without sepsis: 36.76810930704536
varience of temp of patients without sepsis: 0.6712240218361462
the difference of means is:  0.3080081048643919
```

```python
1 # T test
2 S_2 = ((n_sepsis - 1)*sigma_sepsis_temp**2 + (n_non_sepsis - 1)*sigma_not_sepsis_temp**2) / (n_sepsis + n_non_sepsis
3 se_t = np.sqrt(S_2) * np.sqrt((1 / n_sepsis) + (1 / n_non_sepsis))
4 T = d / se_t
5
6 pv = 2 - 2*stats.t.cdf(T, n_sepsis + n_non_sepsis - 2)
7 t_per = stats.t(df=n_sepsis + n_non_sepsis - 2).ppf((0.025, 0.975))
8
9 print('Subsection C\n')
10
11 print('T test:')
12 if  t_per[0] < T < t_per[1]:
13   print('Accept H0')
14 else:
15   print('Reject H0')
16
17
18 print('T is:', T)
19 print('H0 interval:', t_per)
20 print('P-value is:', pv)
21
```

```
Subsection C
```

```
T test:
Reject H0
T is: 16.244895648921773
H0 interval: [-1.96008262  1.96008262]
P-value is: 0.0
```

```python
1
2 sepsis_hr = train_data[train_data['SepsisLabel'] == 1]
3 sepsis_hr = sepsis_hr.HR
4 sepsis_hr = np.array(sepsis_hr)
5
6
7 not_sepsis_hr= train_data[train_data['SepsisLabel'] == 0]
8 not_sepsis_hr = not_sepsis_hr.HR
9 not_sepsis_hr = np.array(not_sepsis_hr)
10
11 mu_sepsis_hr = np.mean(sepsis_hr)
12 mu_non_sepsis_hr = np.mean(not_sepsis_hr)
13
14 n_sepsis = len(sepsis_hr)
15 n_non_sepsis = len(not_sepsis_hr)
16
17 sigma_sepsis_hr = math.sqrt(np.sum((sepsis_hr - mu_sepsis_hr)**2) / n_sepsis)
18 sigma_not_sepsis_hr = math.sqrt(np.sum((not_sepsis_hr - mu_non_sepsis_hr)**2) / n_non_sepsis)
19
20 d = mu_sepsis_hr - mu_non_sepsis_hr
21
22 print('mean of hr of patients with sepsis:', mu_sepsis_hr)
23 print('varience of hr of patients with sepsis:', sigma_sepsis_hr)
24 print('mean of hr of patients without sepsis:', mu_non_sepsis_hr)
25 print('varience of hr of patients without sepsis:', sigma_not_sepsis_hr)
26 print('the difference of means is: ', d)
```

```
mean of hr of patients with sepsis: 90.41795557200774
varience of hr of patients with sepsis: 18.931813129696742
mean of hr of patients without sepsis: 82.79727610503689
varience of hr of patients without sepsis: 17.089776994386146
the difference of means is:  7.620679466970856
```

```python
1 # T test
2 S_2 = ((n_sepsis - 1)*sigma_sepsis_hr**2 + (n_non_sepsis - 1)*sigma_not_sepsis_hr**2) / (n_sepsis + n_non_sepsis - 2)
3 se_t = np.sqrt(S_2) * np.sqrt((1 / n_sepsis) + (1 / n_non_sepsis))
4 T = d / se_t
5
6 pv = 2 - 2*stats.t.cdf(T, n_sepsis + n_non_sepsis - 2)
7 t_per = stats.t(df=n_sepsis + n_non_sepsis - 2).ppf((0.025, 0.975))
8
9 print('Subsection C\n')
10
11 print('T test:')
12 if  t_per[0] < T < t_per[1]:
13   print('Accept H0')
14 else:
15   print('Reject H0')
16
17
18 print('T is:', T)
19 print('H0 interval:', t_per)
20 print('P-value is:', pv)
21
```

```
Subsection C

T test:
Reject H0
T is: 16.041363973288988
H0 interval: [-1.96008262  1.96008262]
P-value is: 0.0
```

```python
1 import pandas as pd
2 df = pd.DataFrame({'a':[1,1,0,1,0], 'b':[1,0,1,1,1]})
3 df['c'] = df['a'] | df['b']
4
```

```
Index(['a', 'b', 'c'], dtype='object')
```