

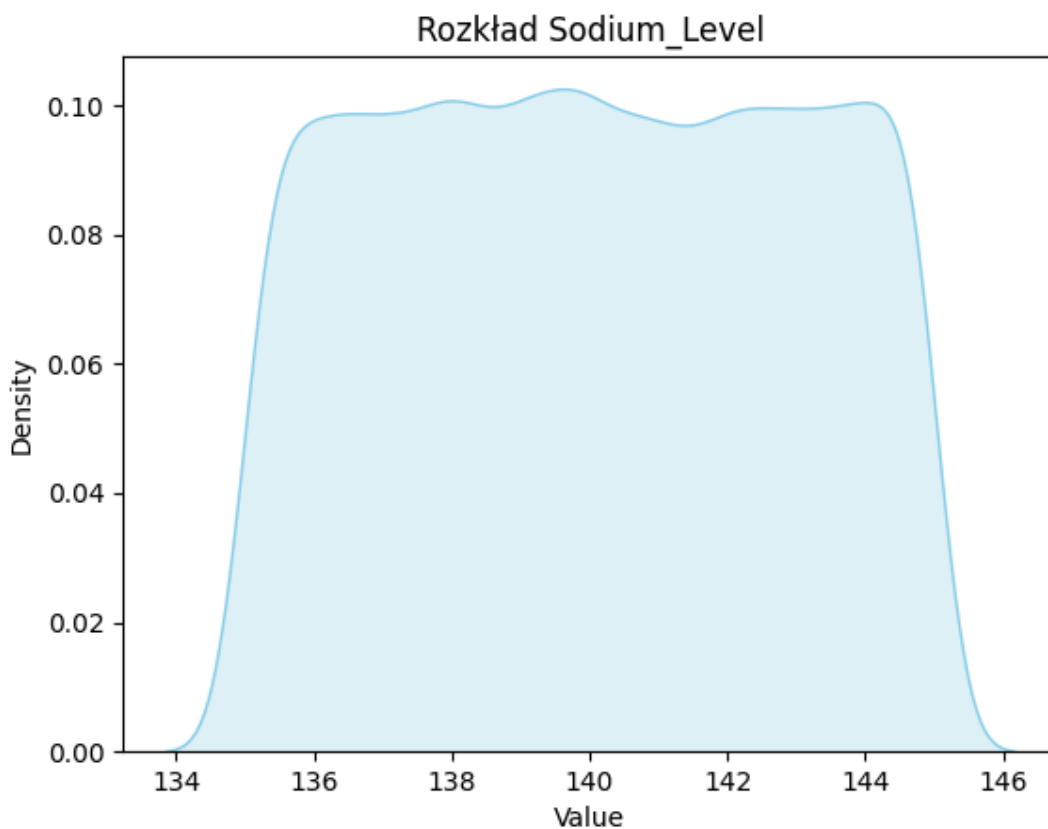
Prototypowanie

Analiza danych

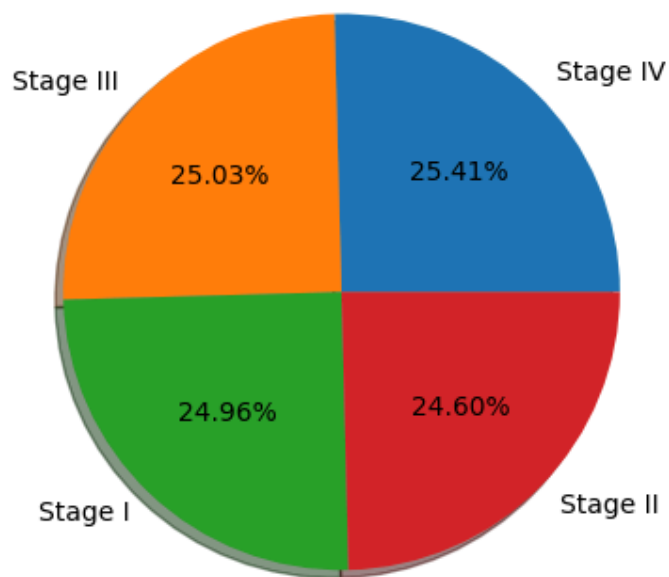
Dokładne efekty działania można zobaczyć w `Analizing_data_set.ipynb`.

Własna analiza danych

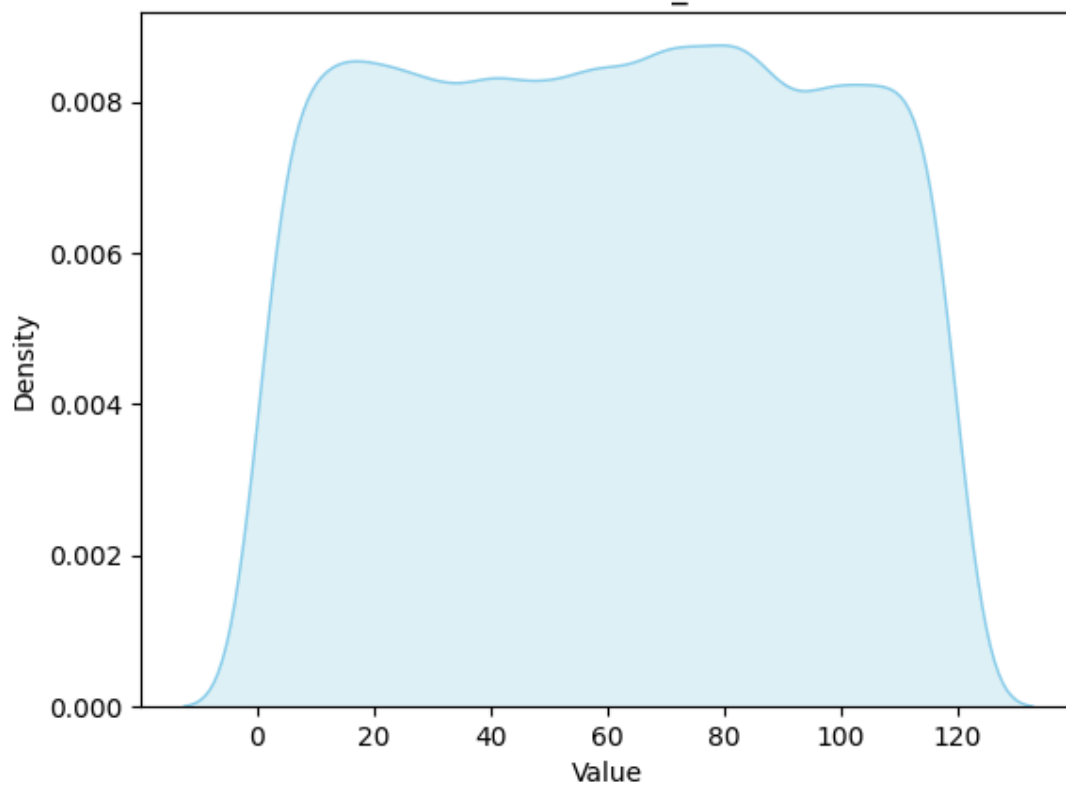
Jak można zobaczyć trzeba było zmienić kolumnę `Performance_Status` z numerycznej na kategoryjną ponieważ ma konkretnie 5 różnych wartości (kategorii) z których można wybrać wartość. Dalej można zobaczyć rozkład wartości każdej kolumny które zostały przedstawione w formie wykresów. Można na nich zauważyć, że dane są dość równomiernie rozłożone.



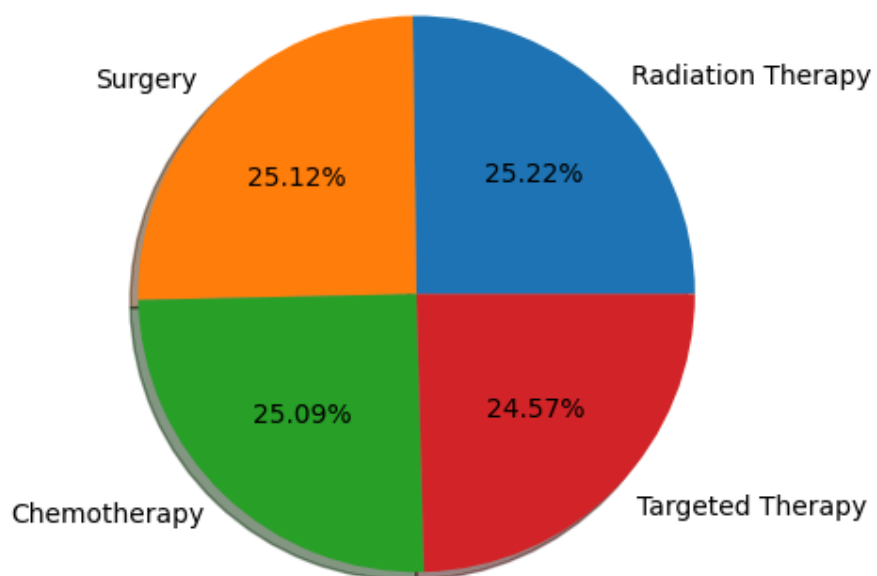
Rozkład Stage



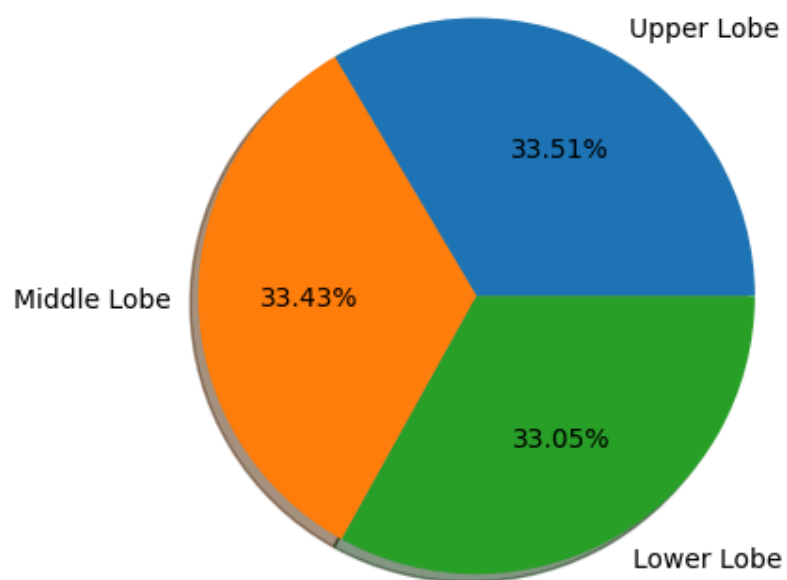
Rozkład Survival_Months

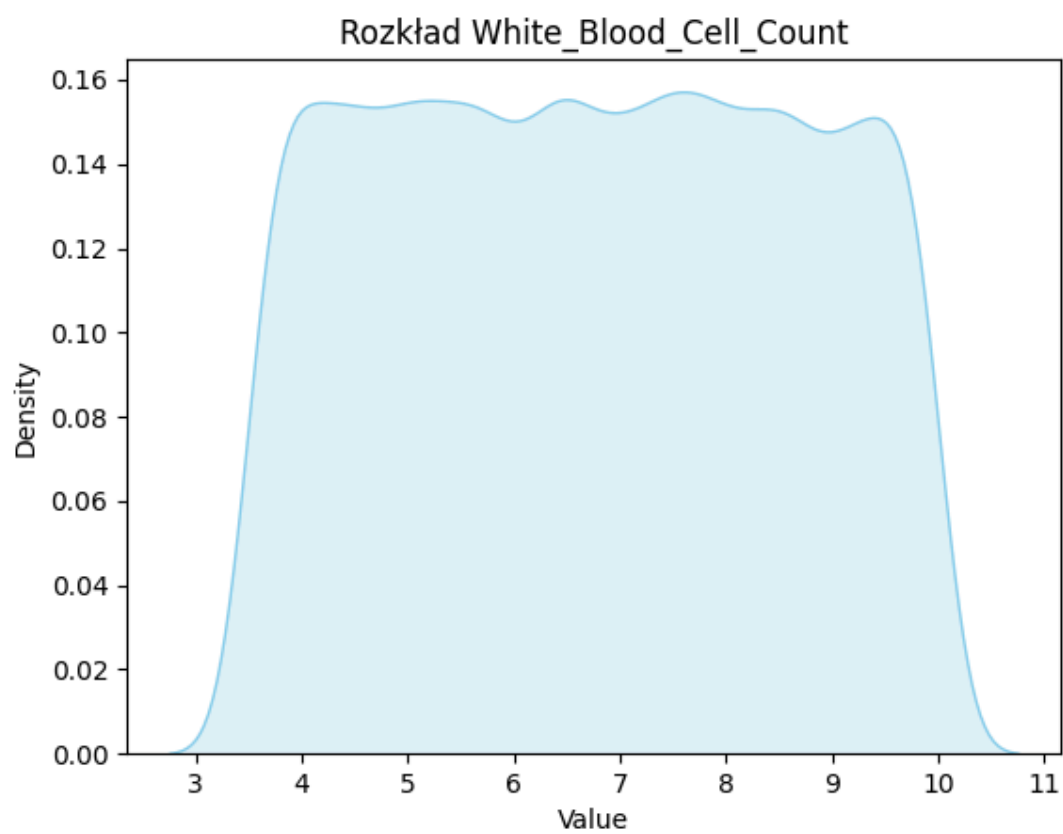
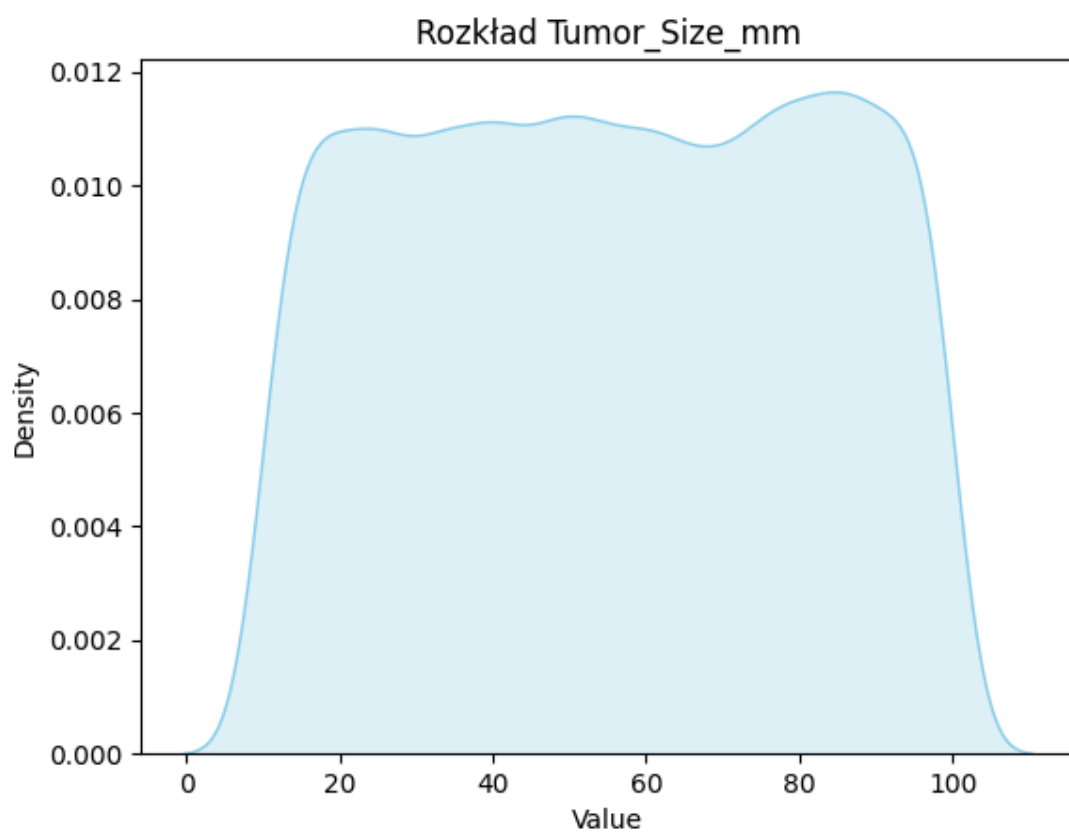


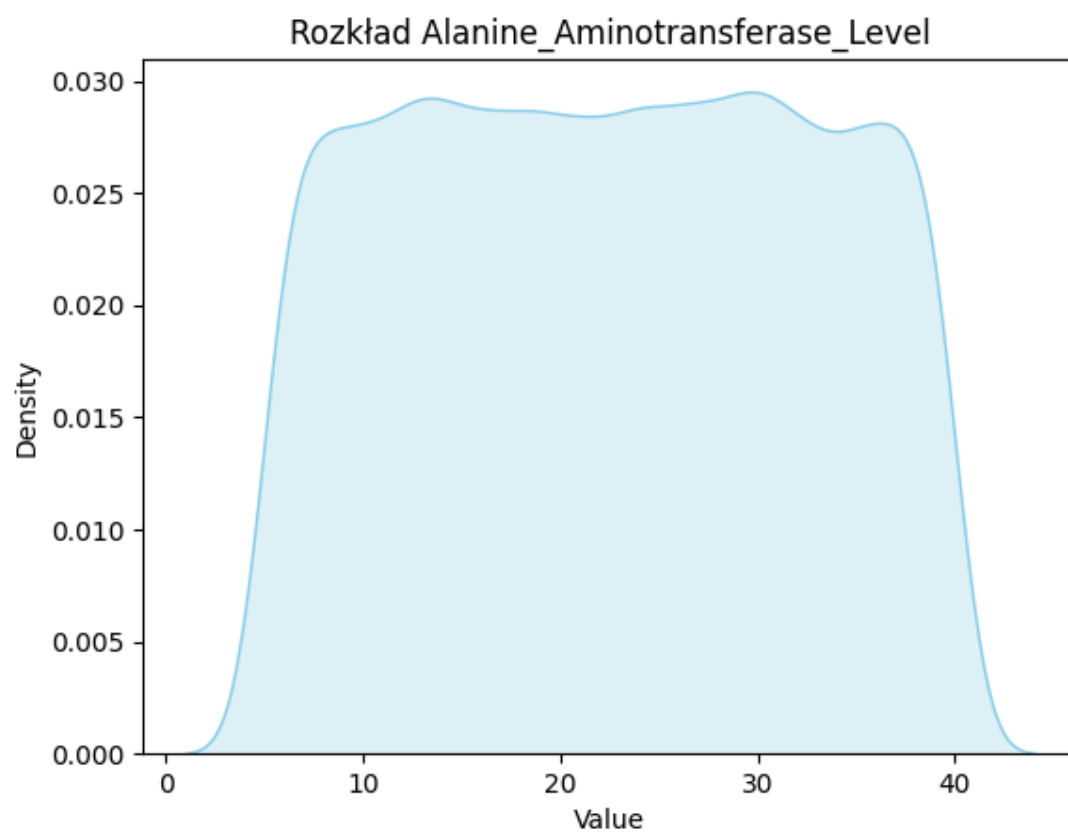
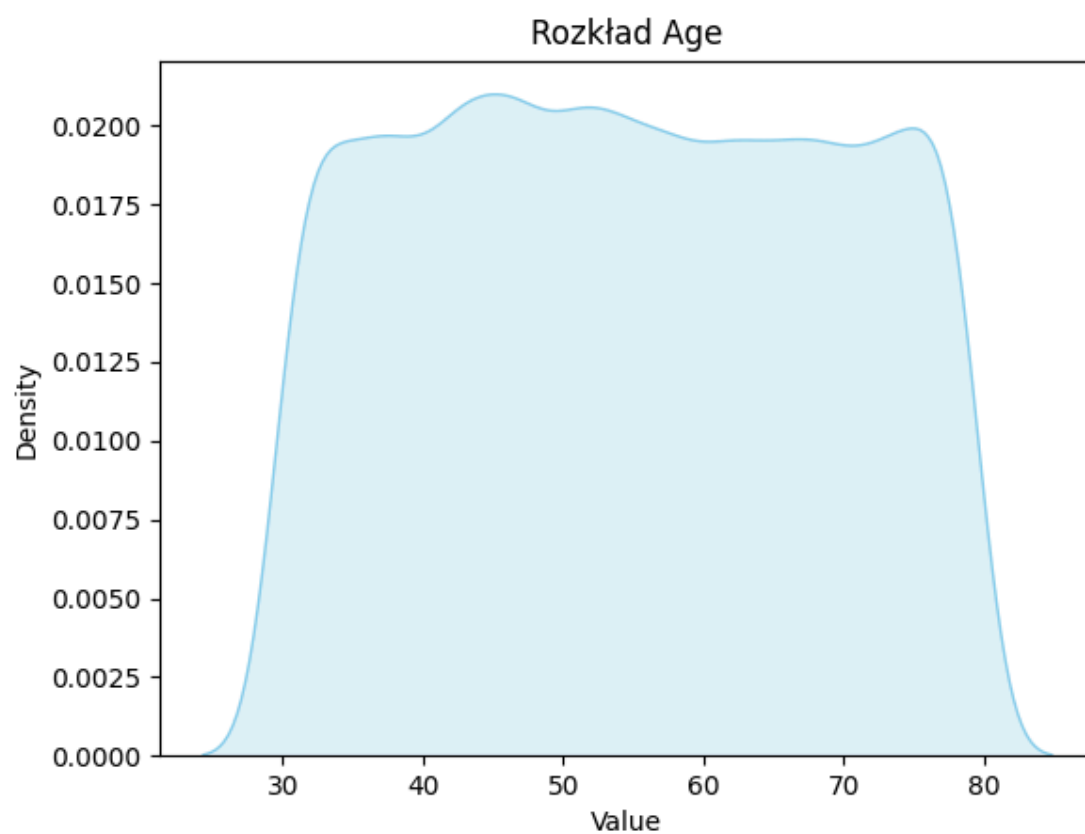
Rozkład Treatment

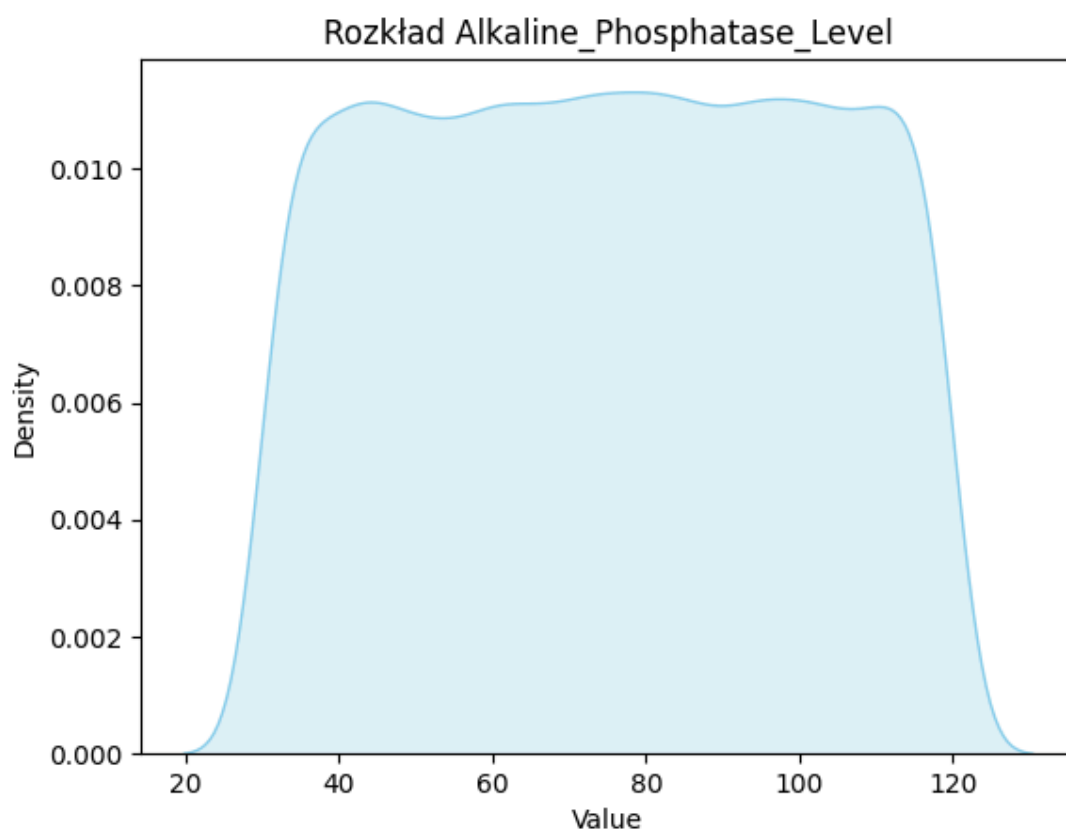
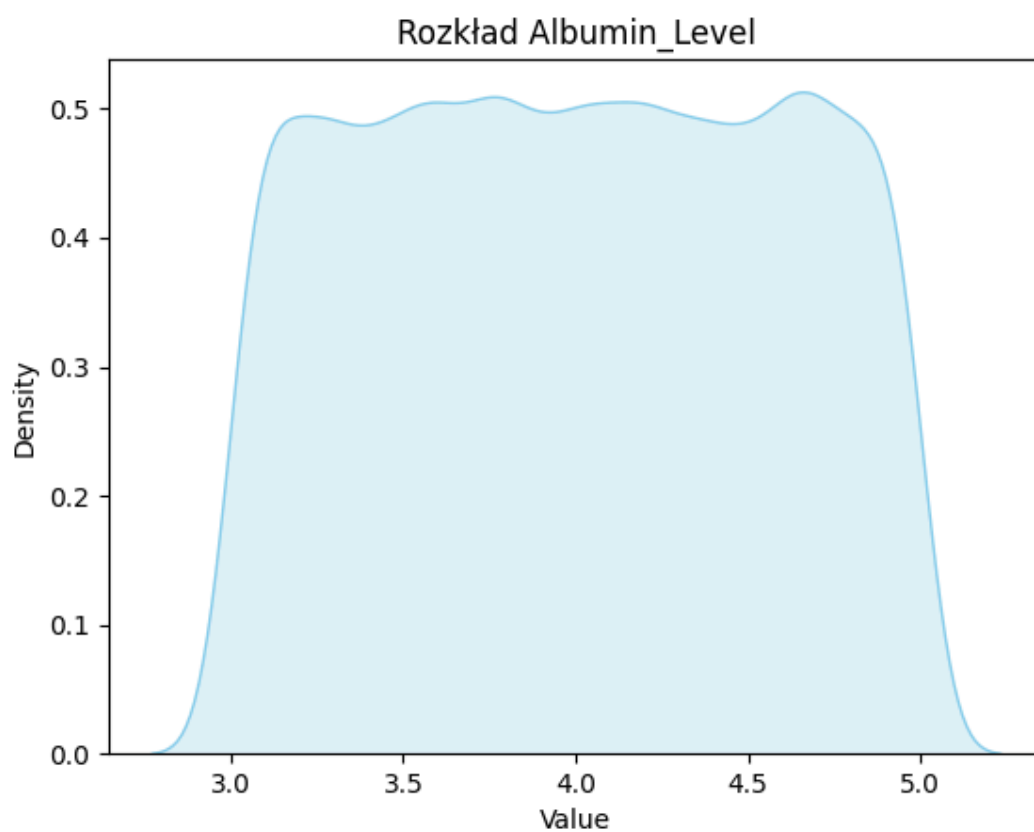


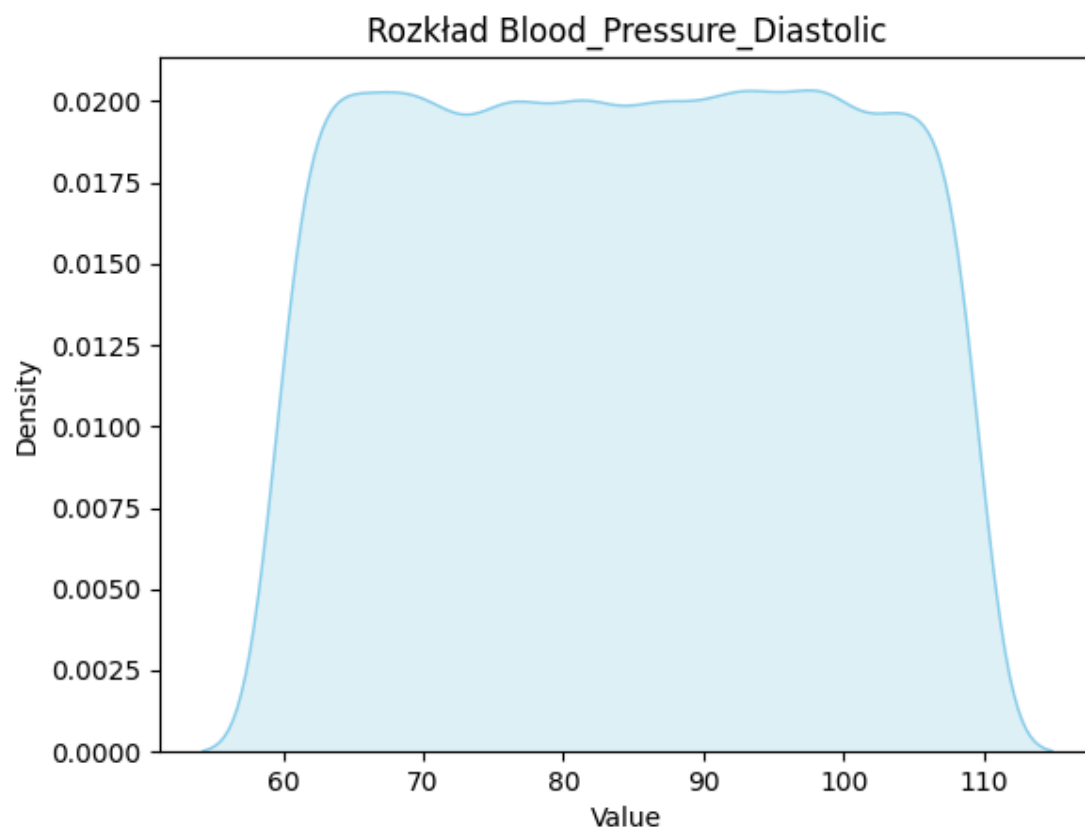
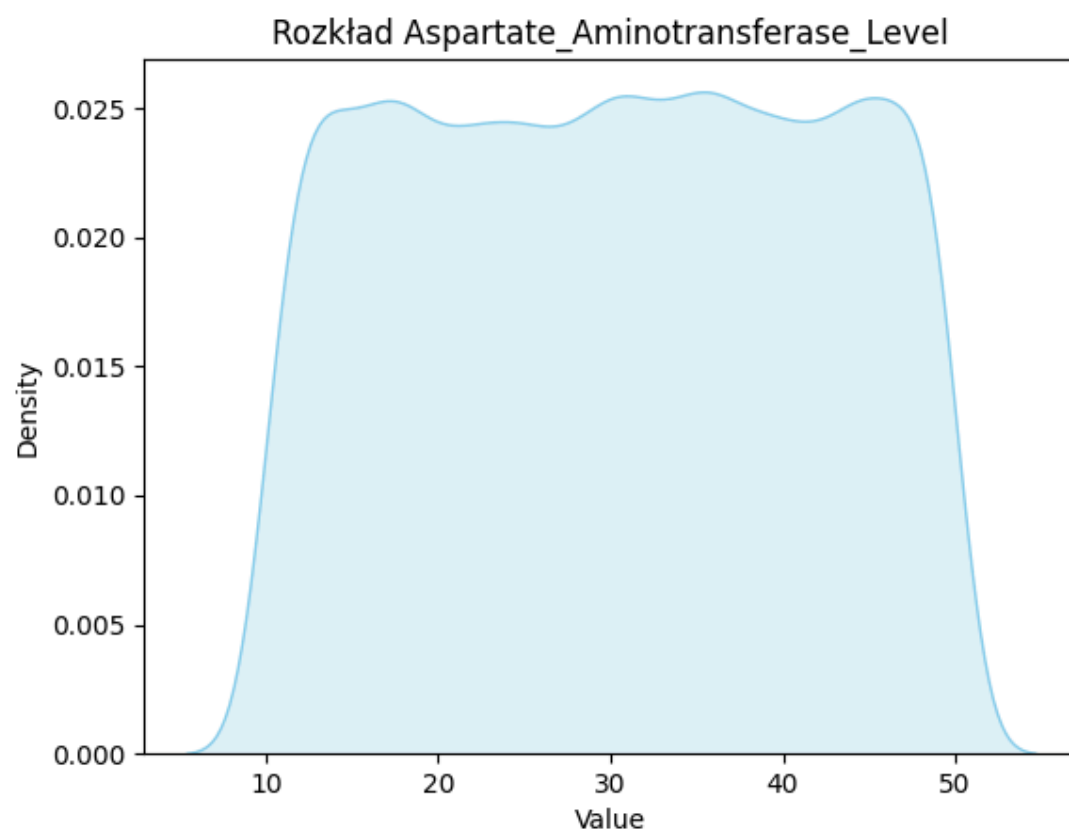
Rozkład Tumor_Location

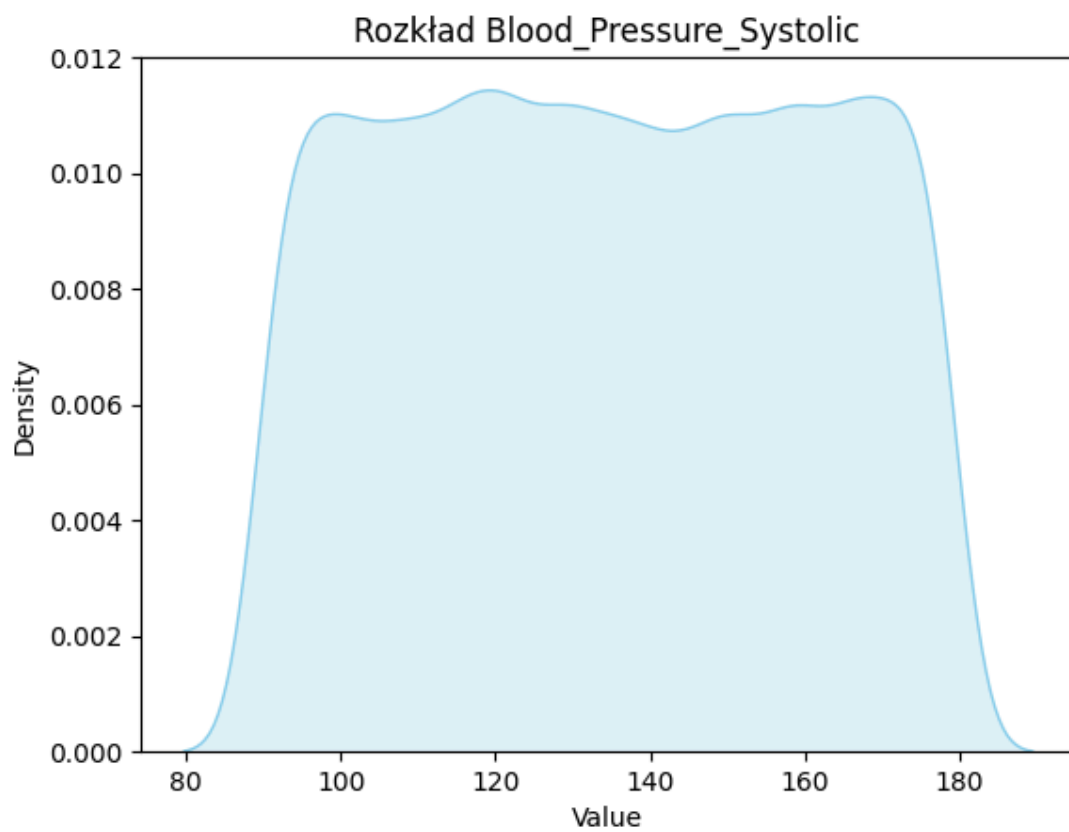
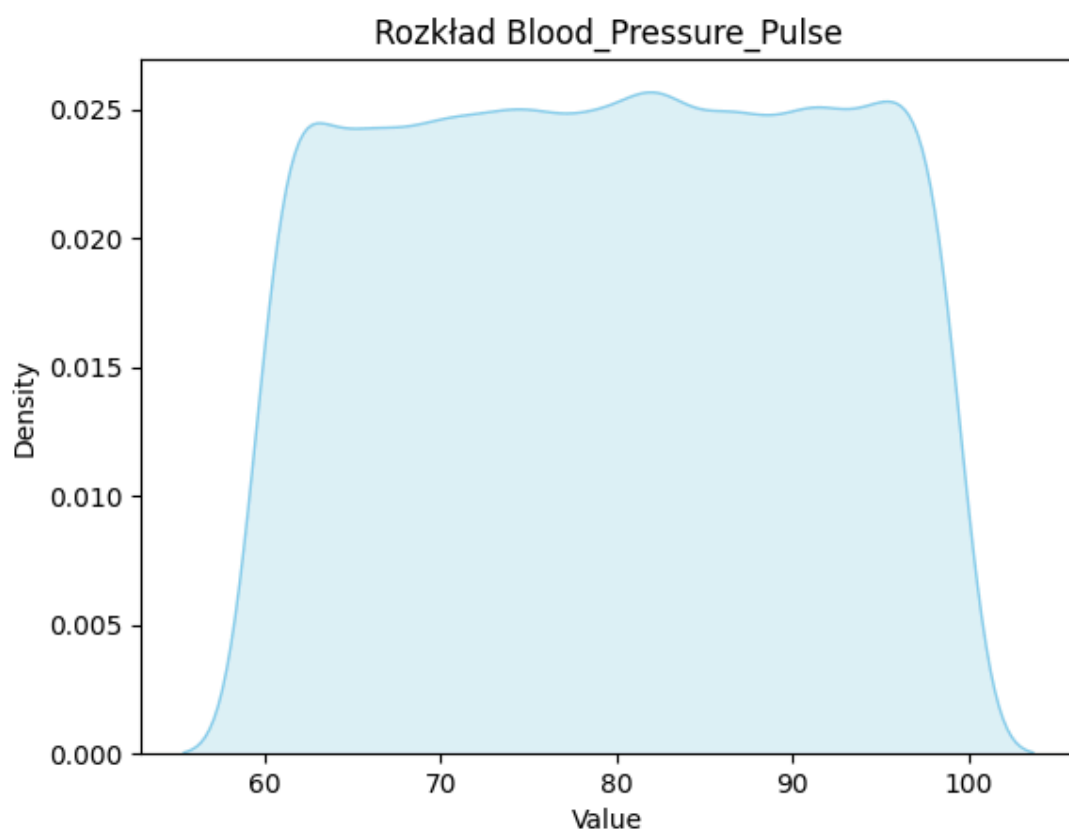


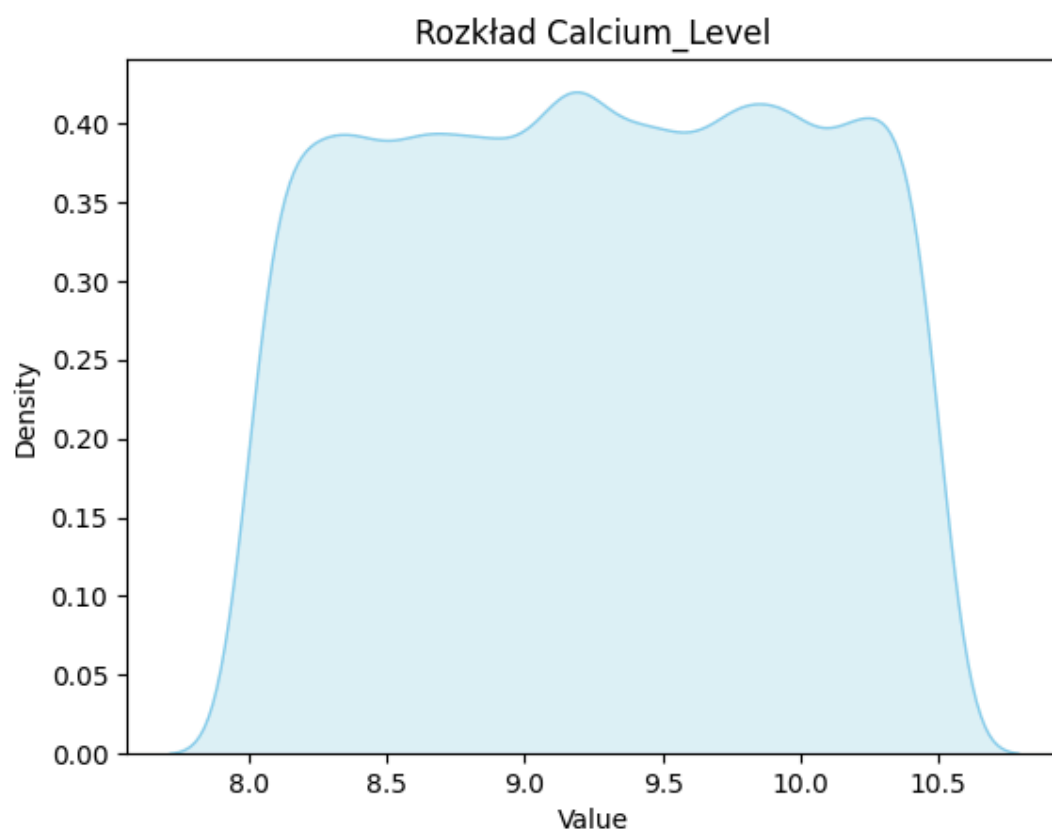




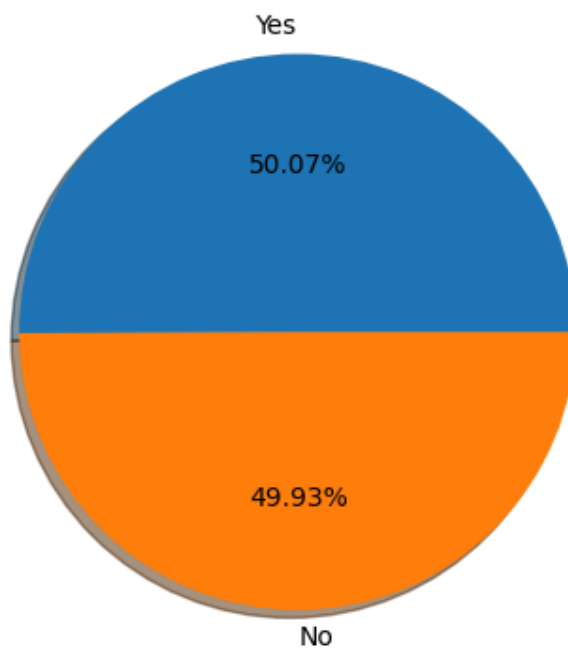




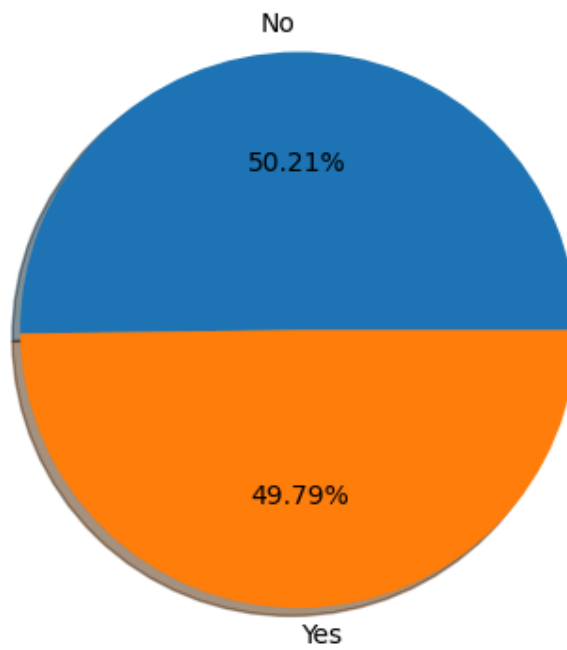




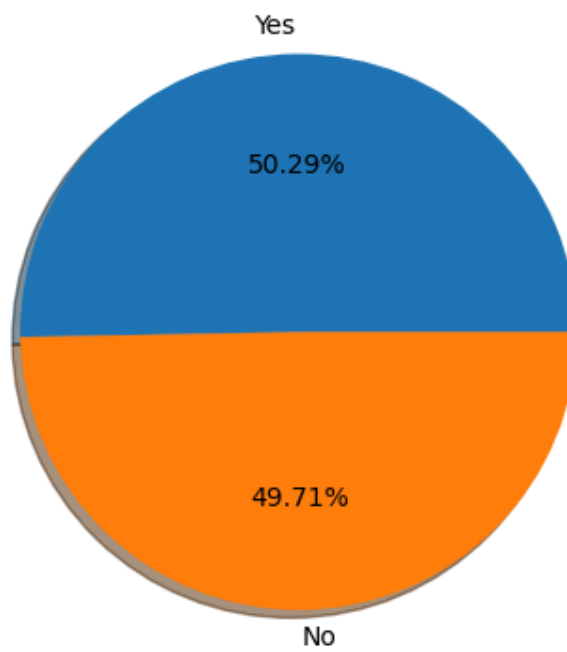
Rozkład Comorbidity_Autoimmune_Disease



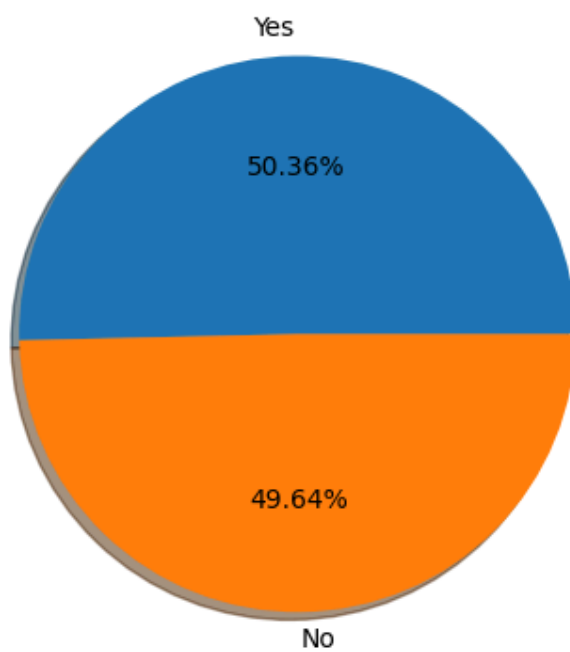
Rozkład Comorbidity_Chronic_Lung_Disease



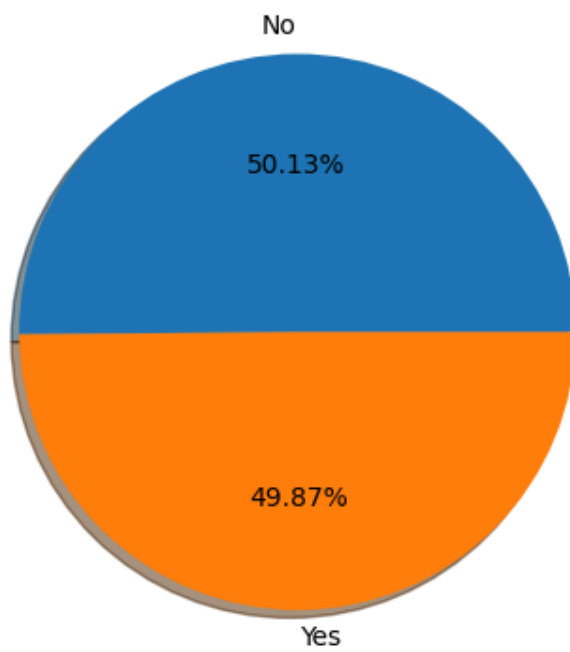
Rozkład Comorbidity_Diabetes



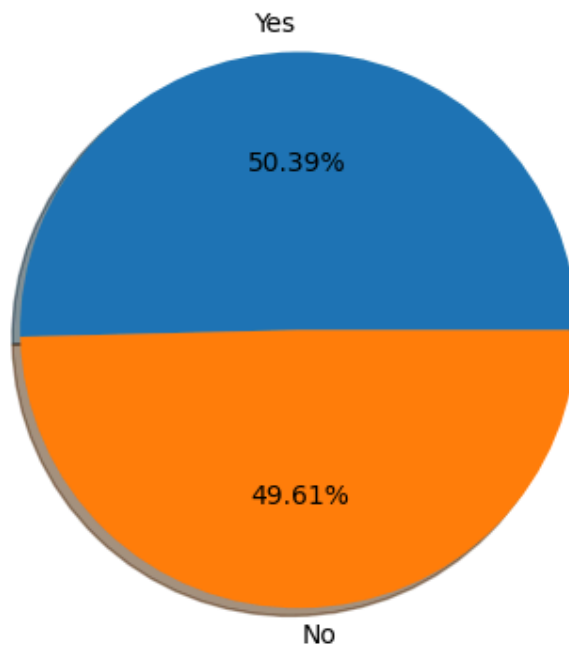
Rozkład Comorbidity_Heart_Disease



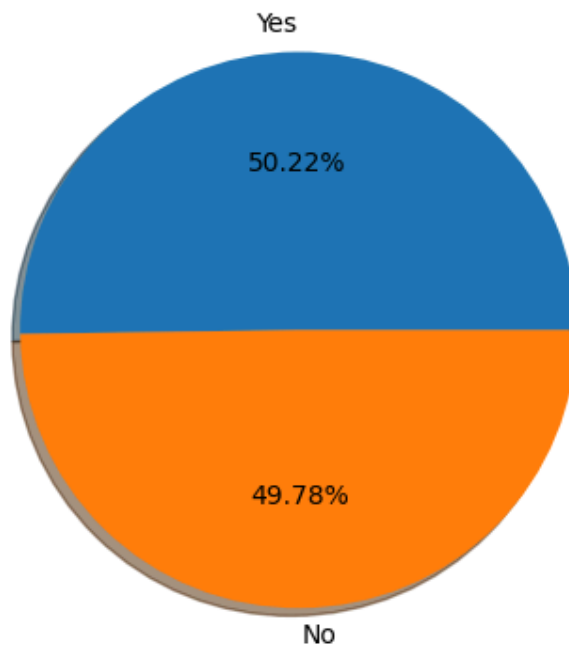
Rozkład Comorbidity_Hypertension

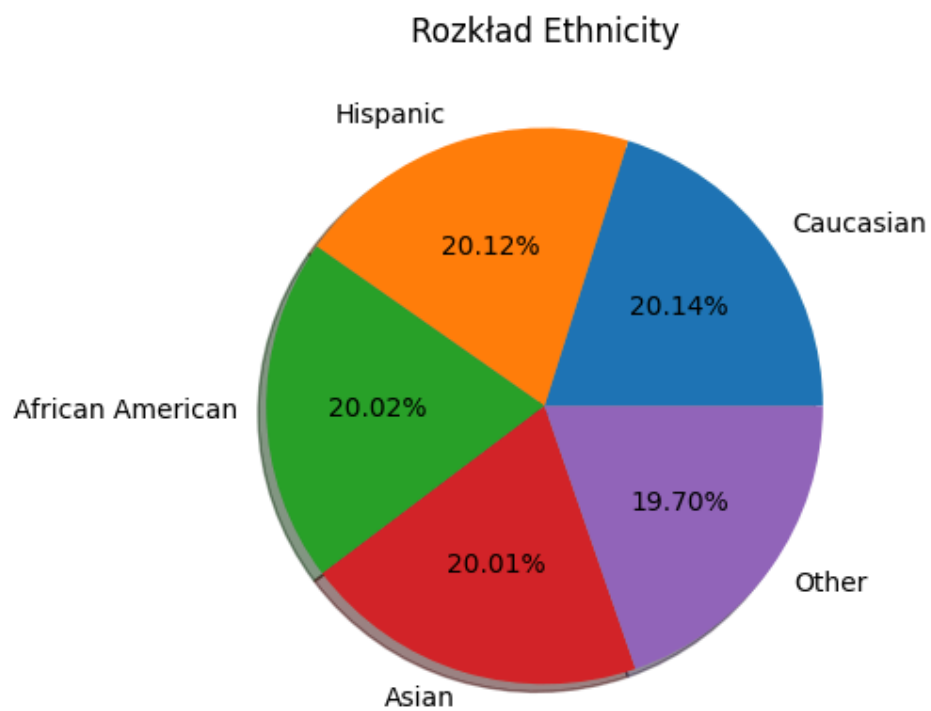
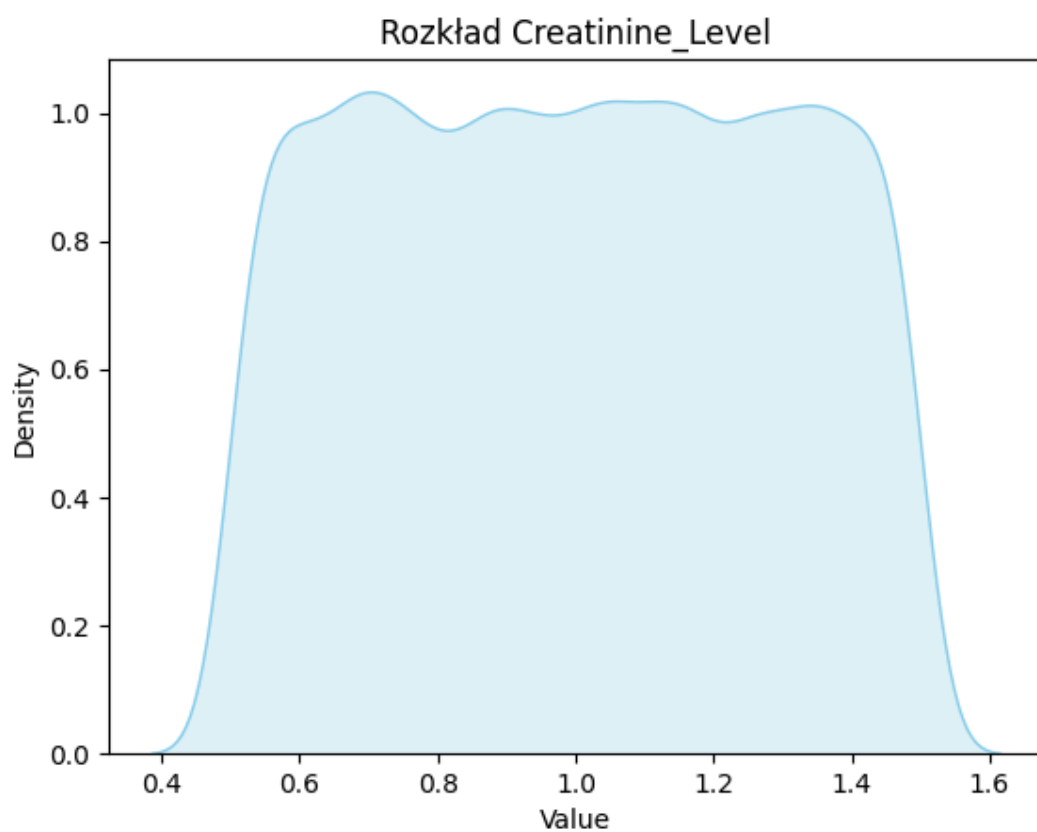


Rozkład Comorbidity_Kidney_Disease

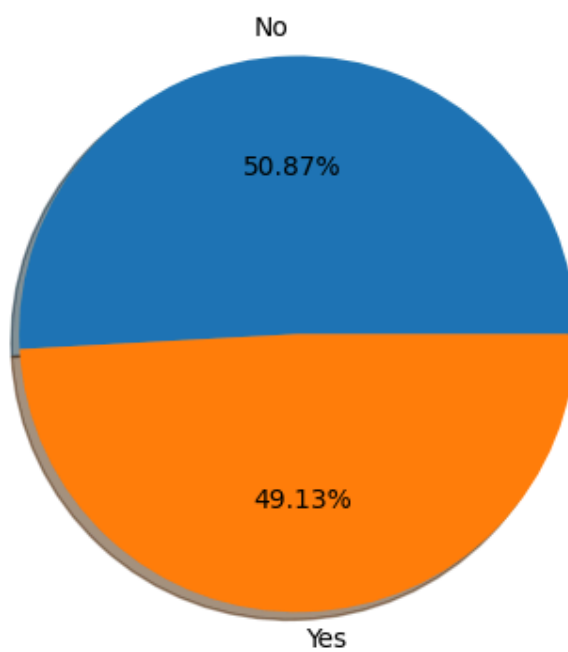


Rozkład Comorbidity_Other

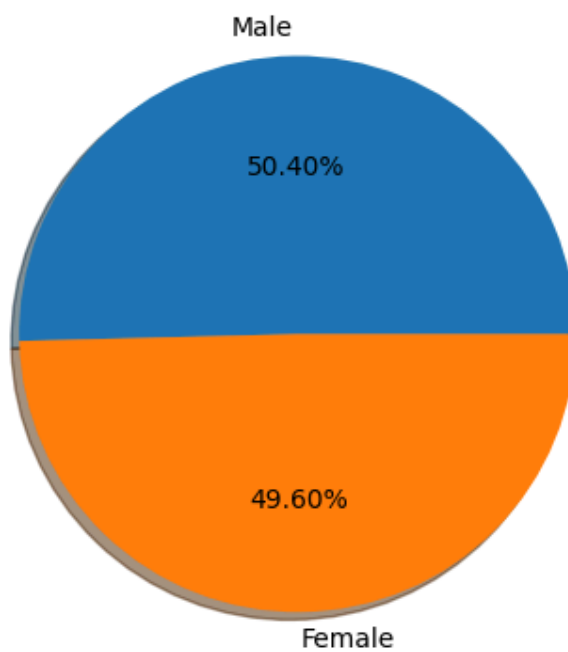


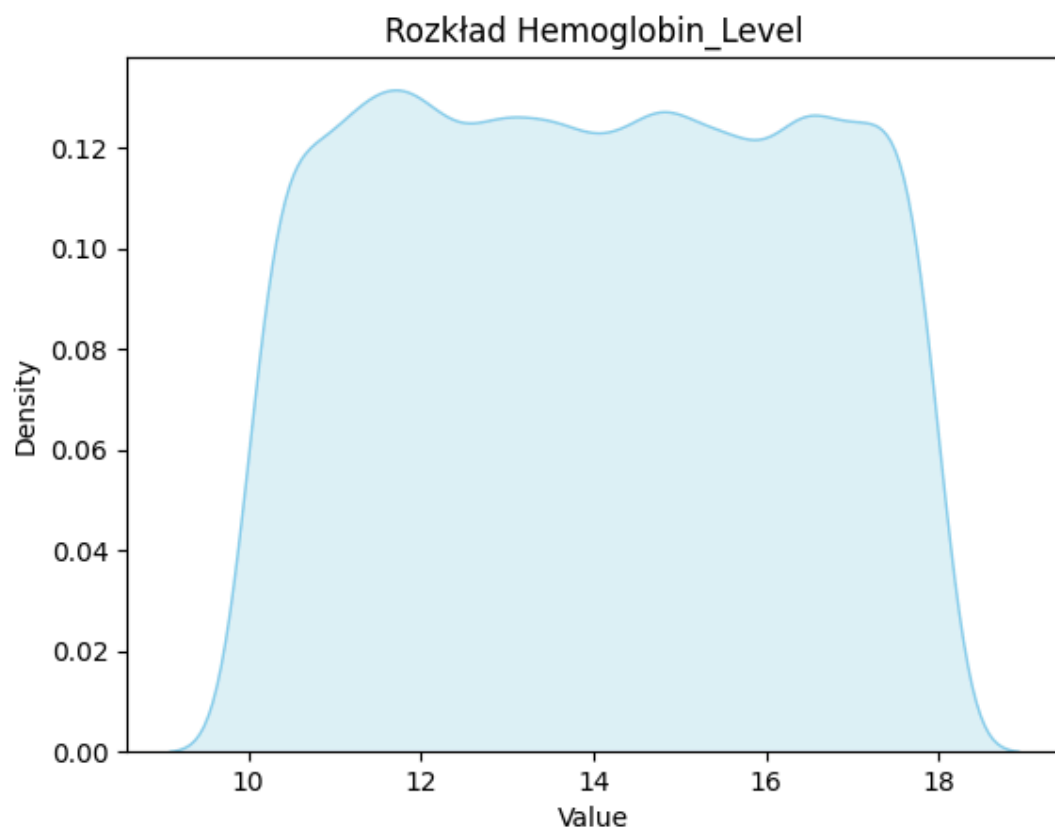
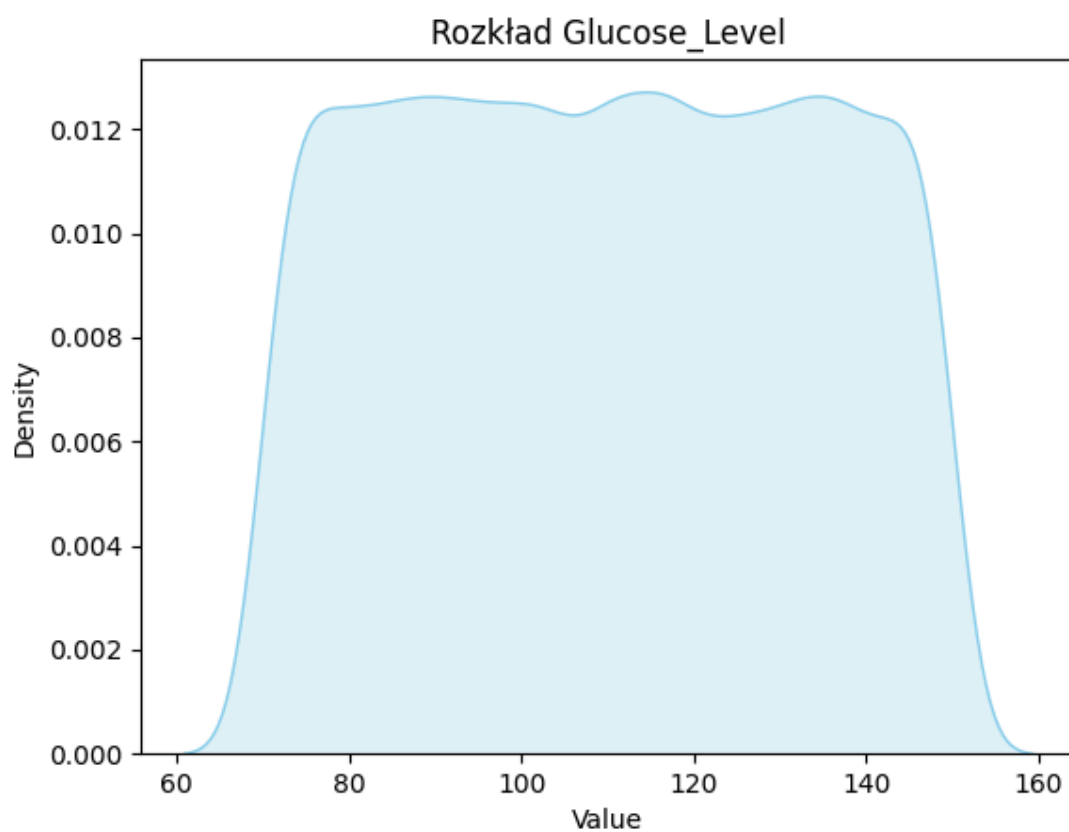


Rozkład Family_History

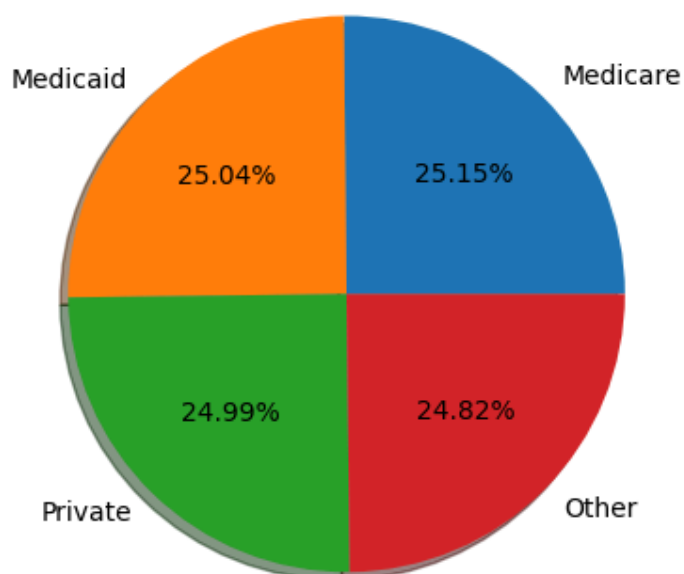


Rozkład Gender

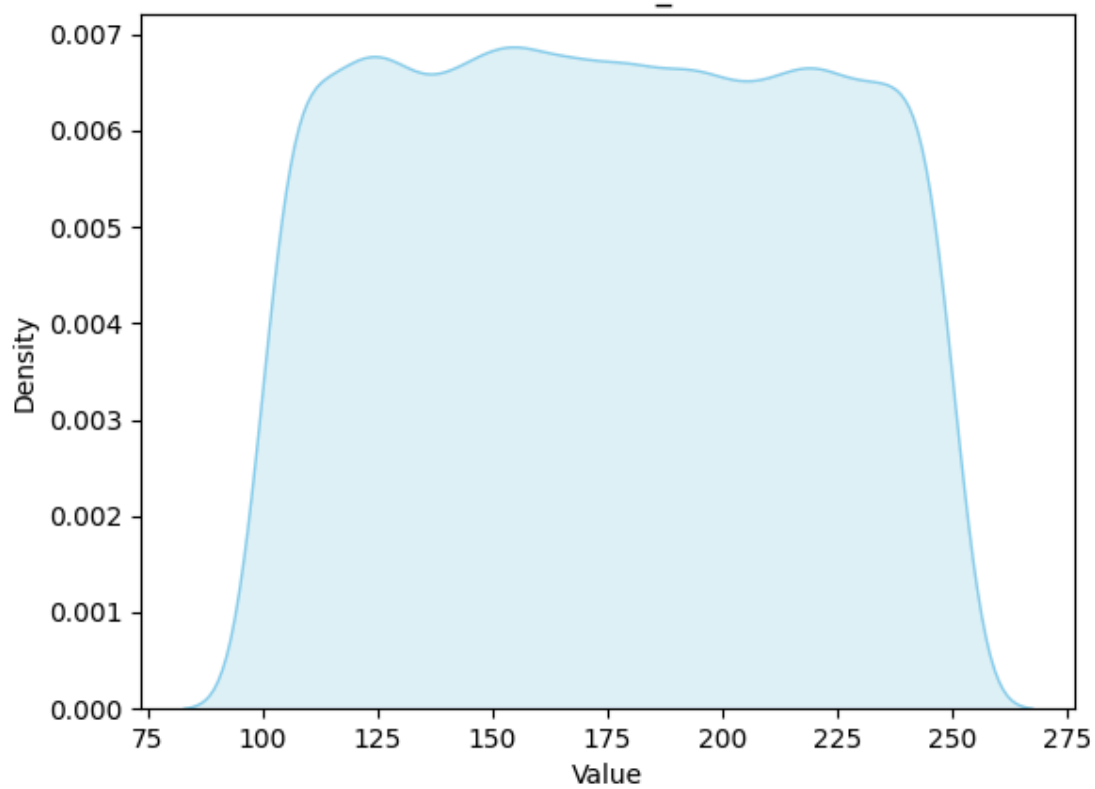




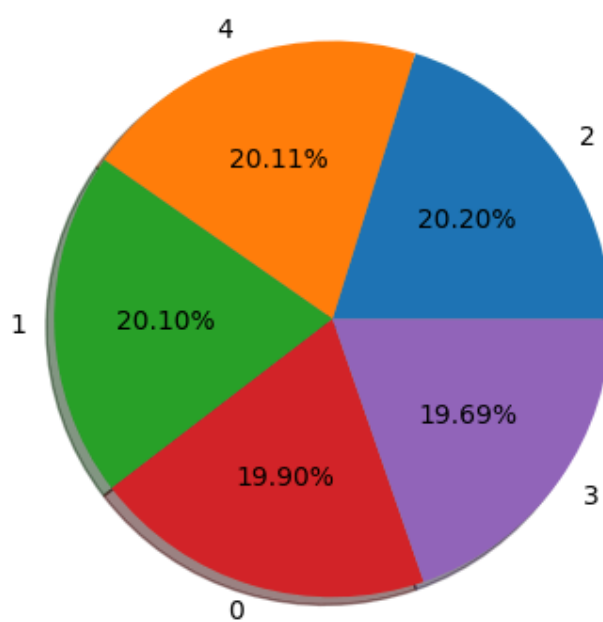
Rozkład Insurance_Type



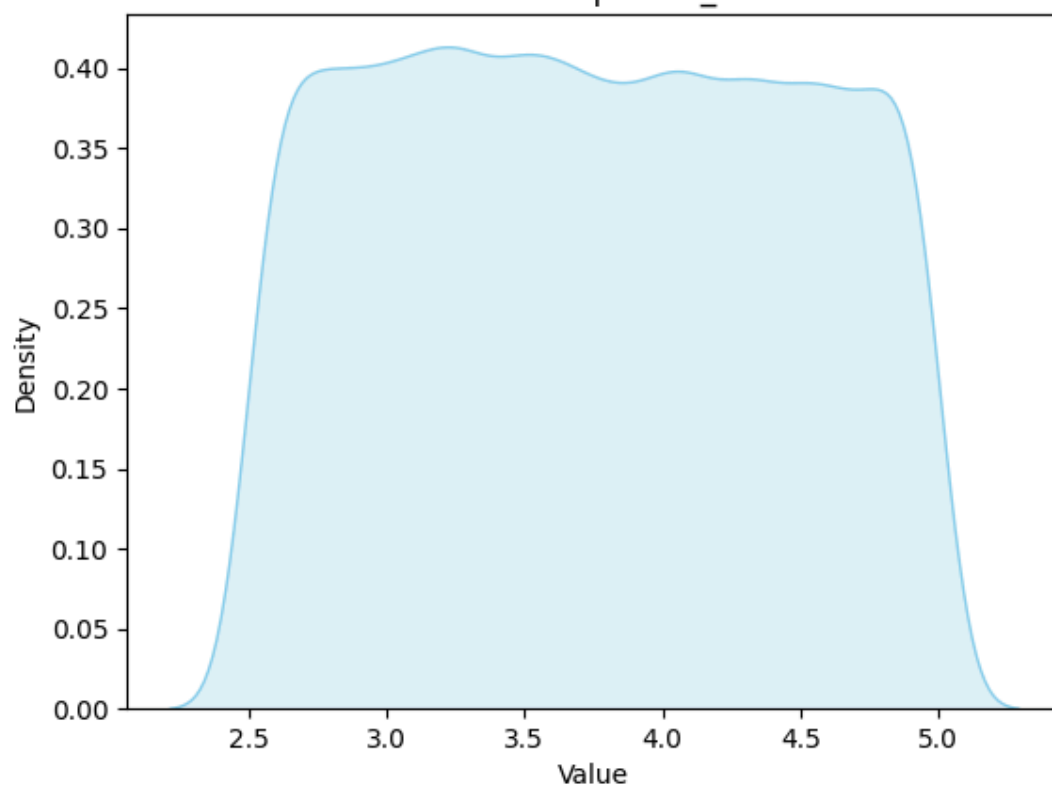
Rozkład LDH_Level

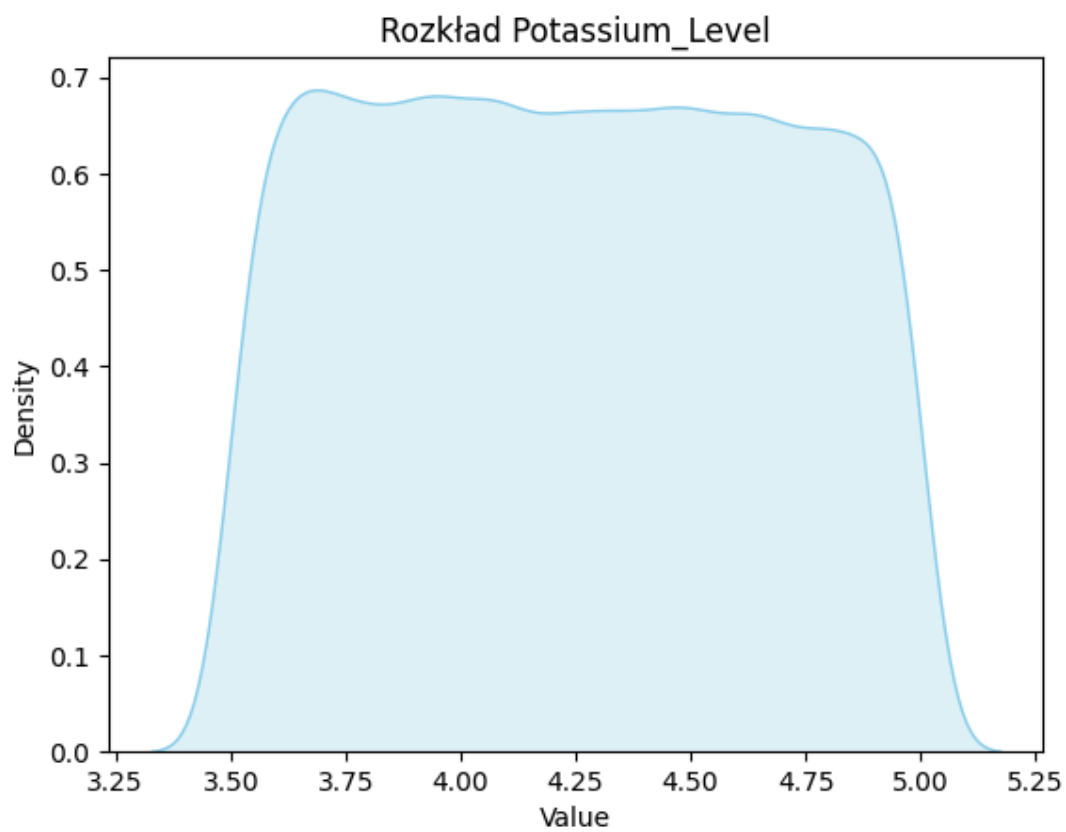
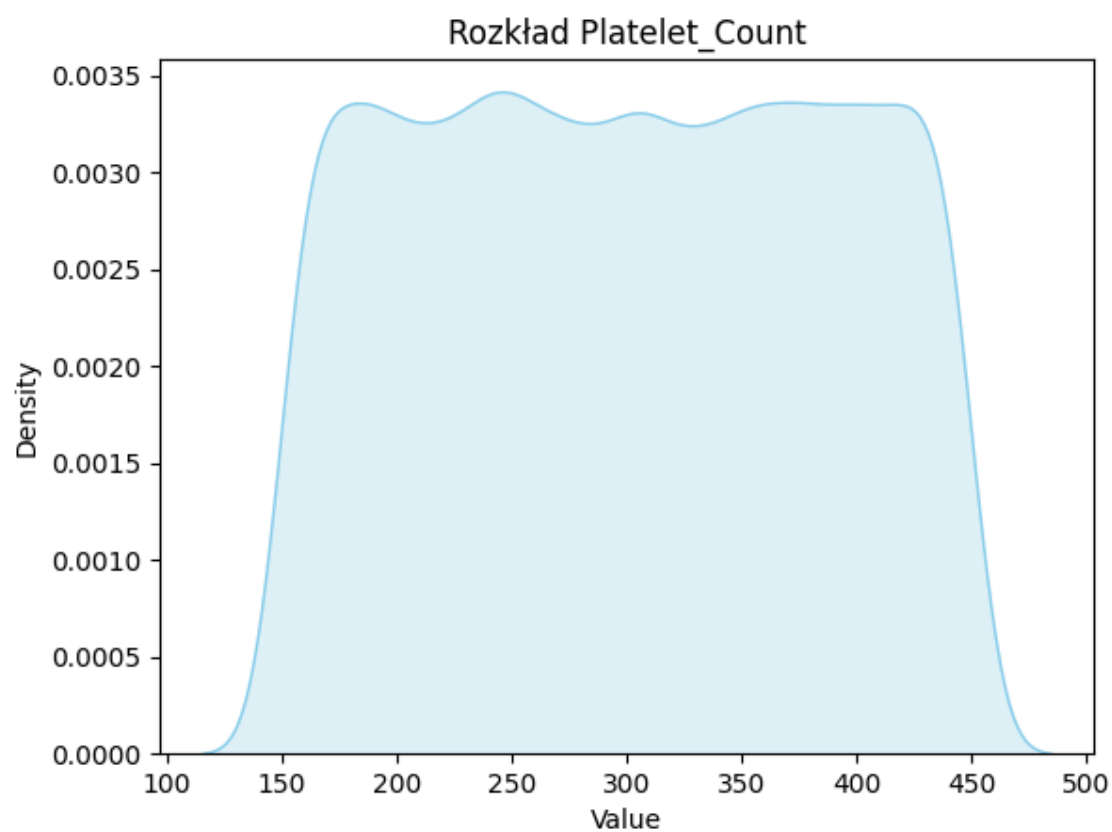


Rozkład Performance_Status

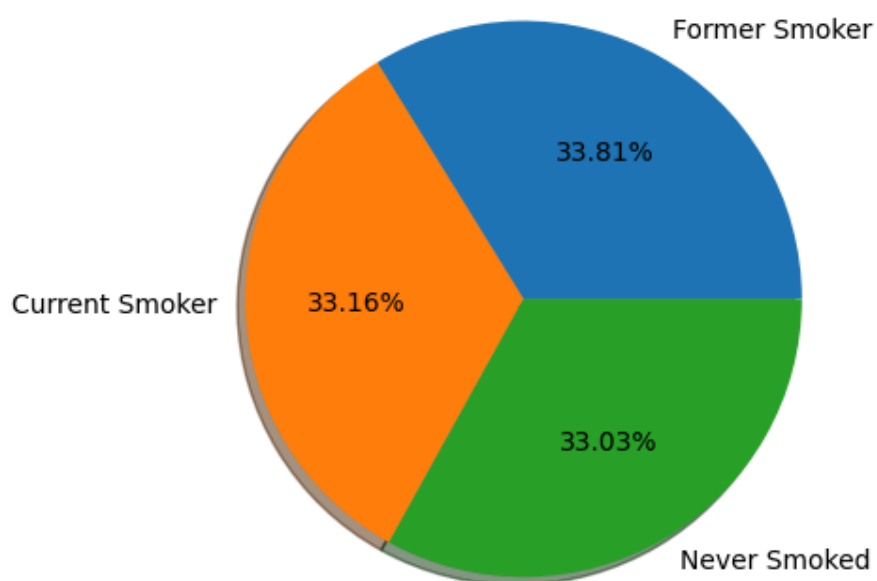


Rozkład Phosphorus_Level

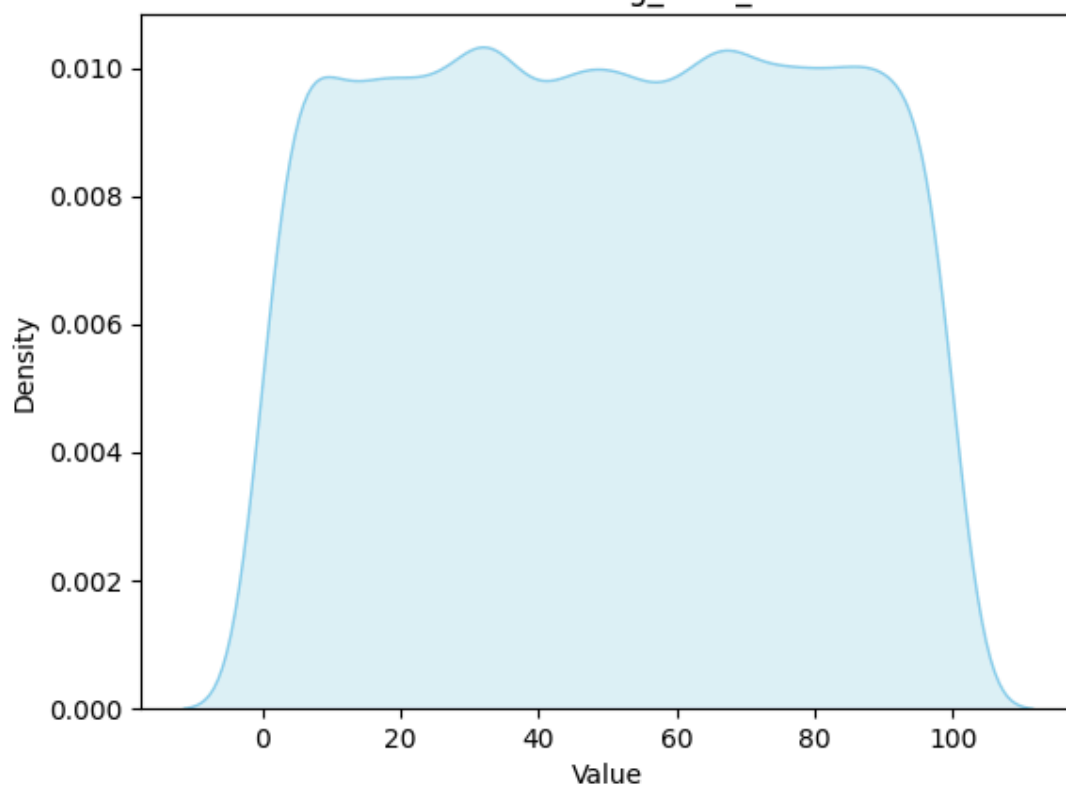




Rozkład Smoking_History

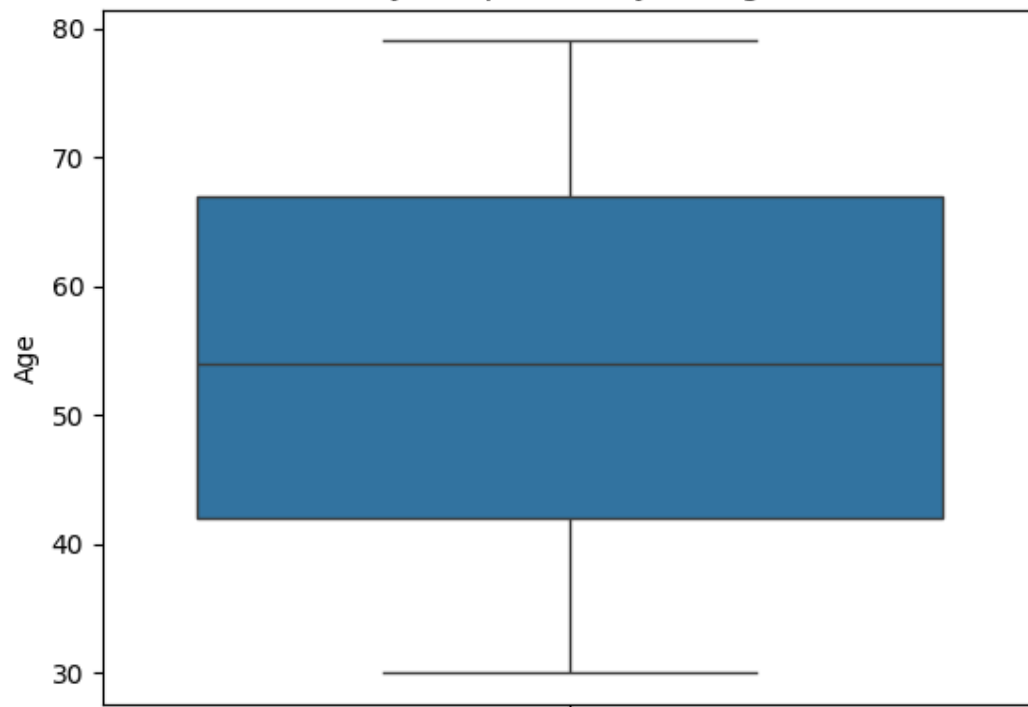


Rozkład Smoking_Pack_Years

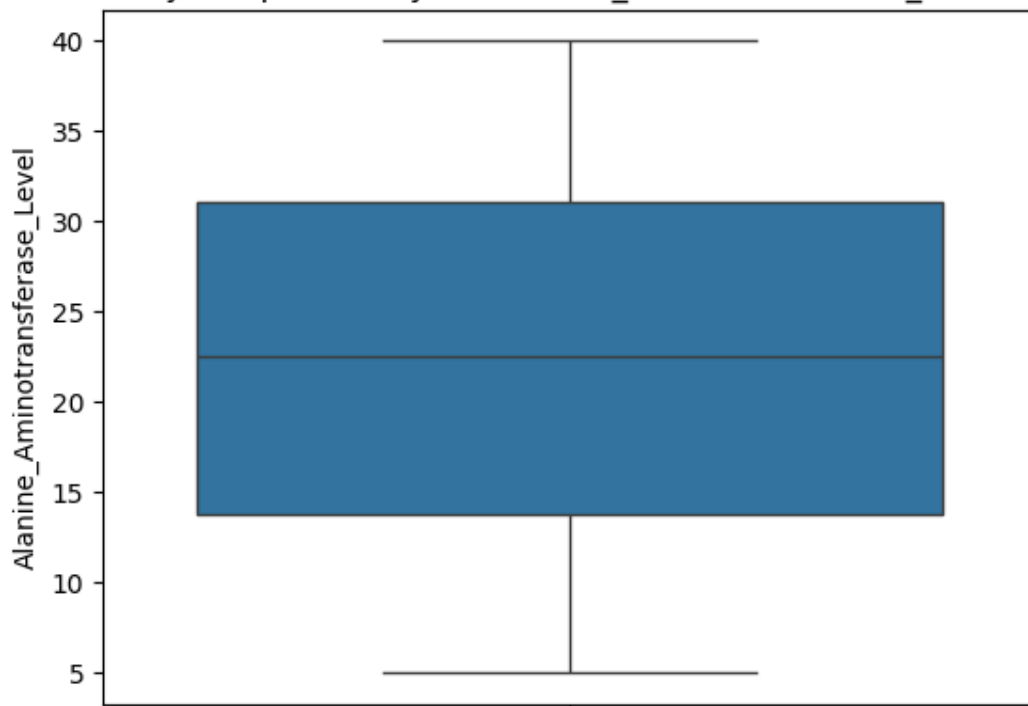


W dalszej części można też zauważyć wykresy pudełkowe dla wartości numerycznych które wydają się podobne do siebie. Można na nich zobaczyć, że minimalna wartość, wartość w $1/4$, wartość w połowie, wartość w $3/4$ i wartość maksymalna są odpowiednio w podobnych odległościach.

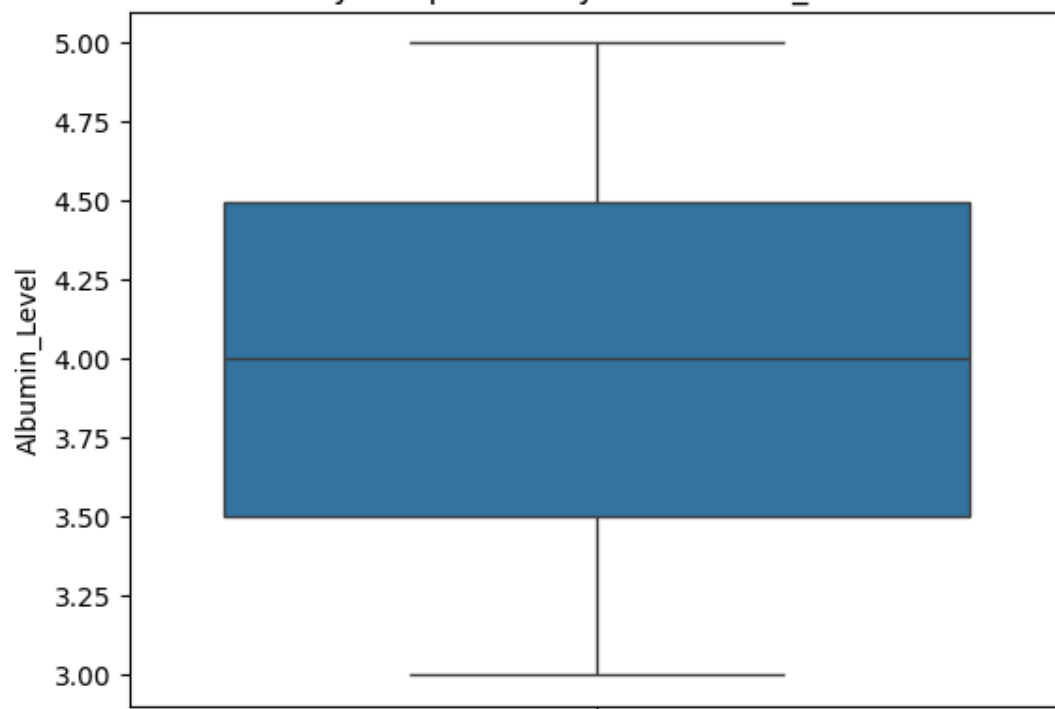
Wykres pudełkowy dla Age



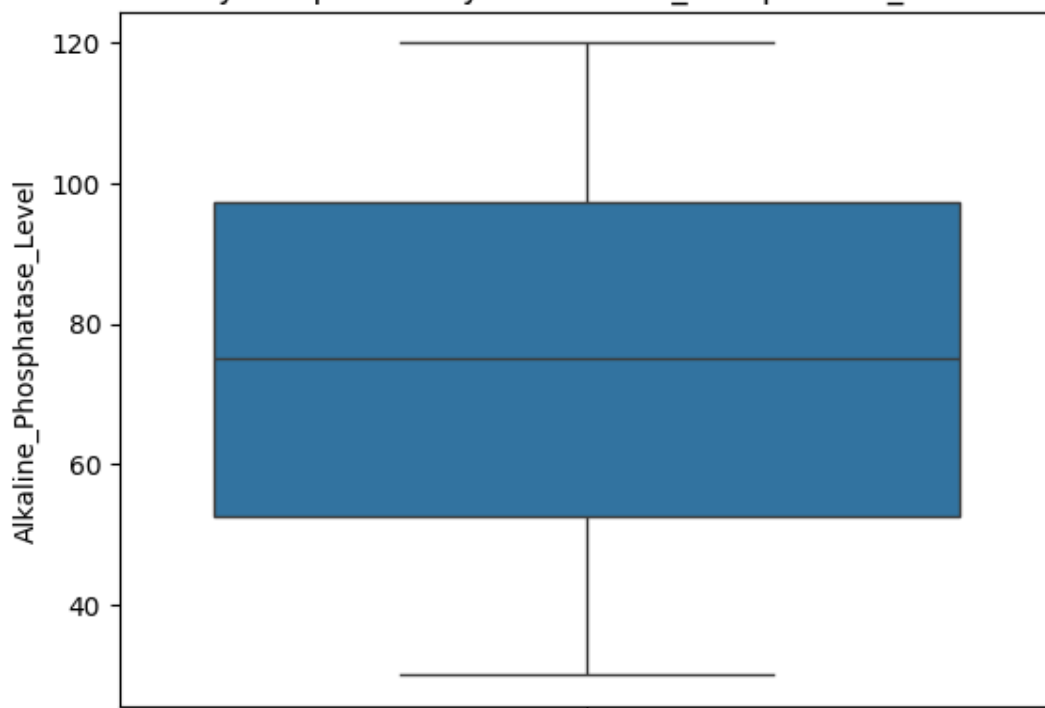
Wykres pudełkowy dla Alanine_Aminotransferase_Level



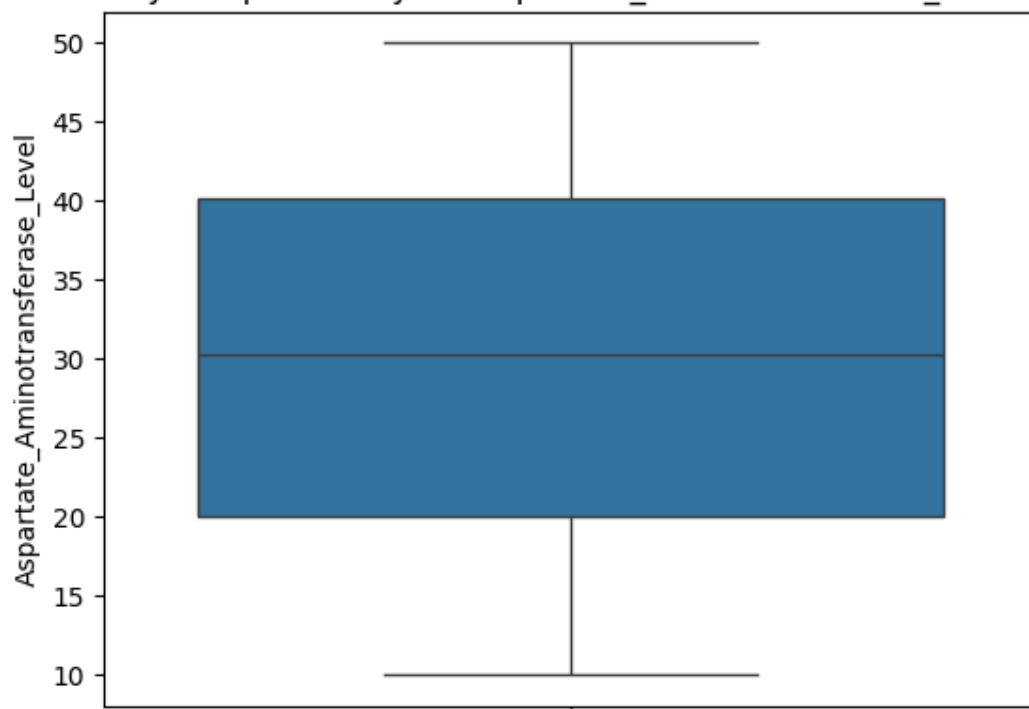
Wykres pudełkowy dla Albumin_Level



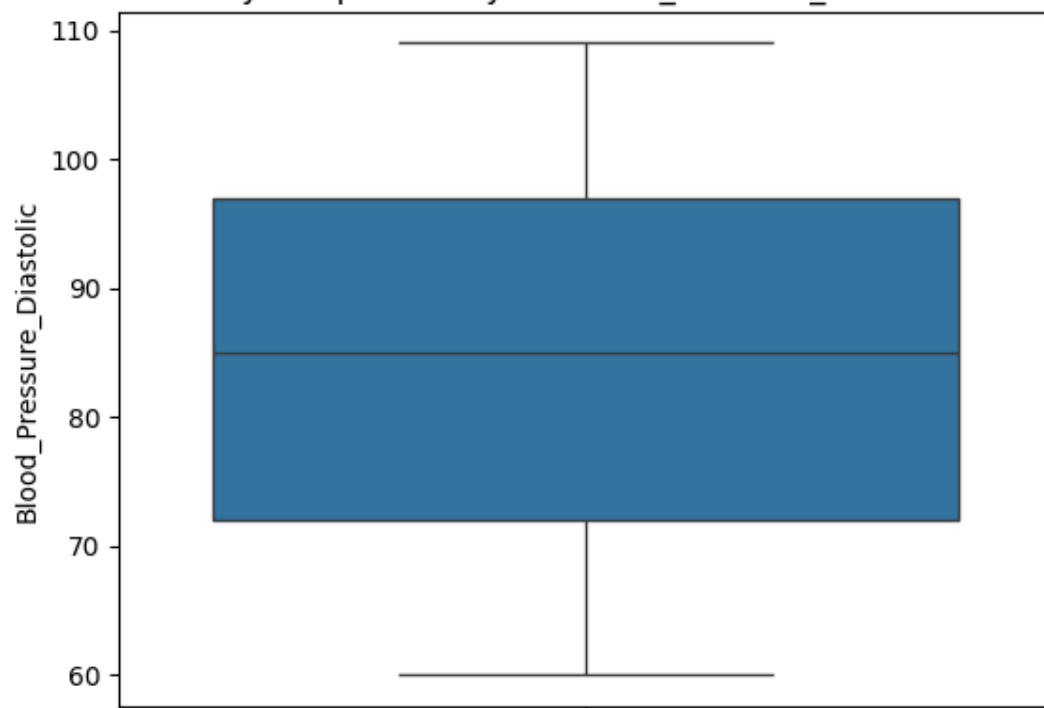
Wykres pudełkowy dla Alkaline_Phosphatase_Level



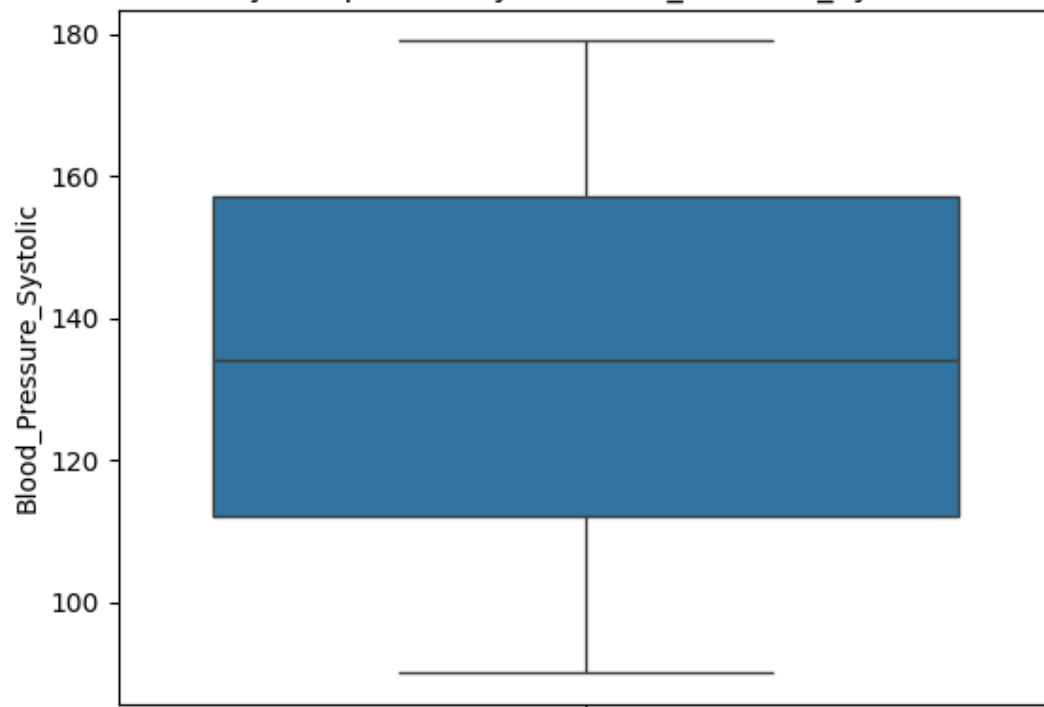
Wykres pudełkowy dla Aspartate_Aminotransferase_Level



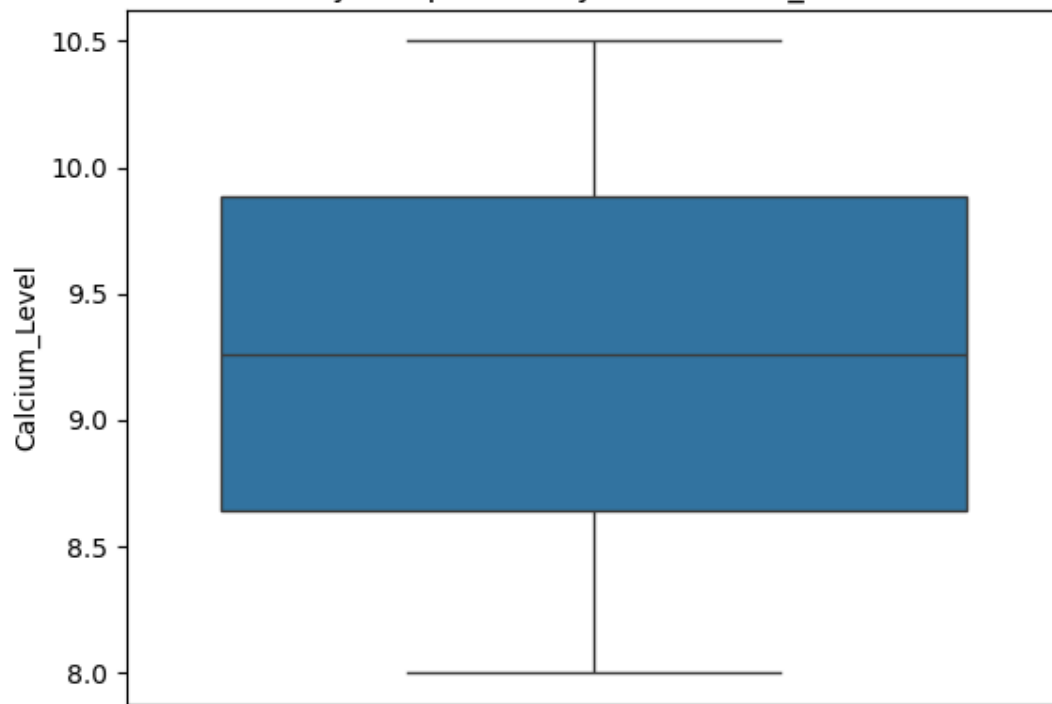
Wykres pudełkowy dla Blood_Pressure_Diastolic



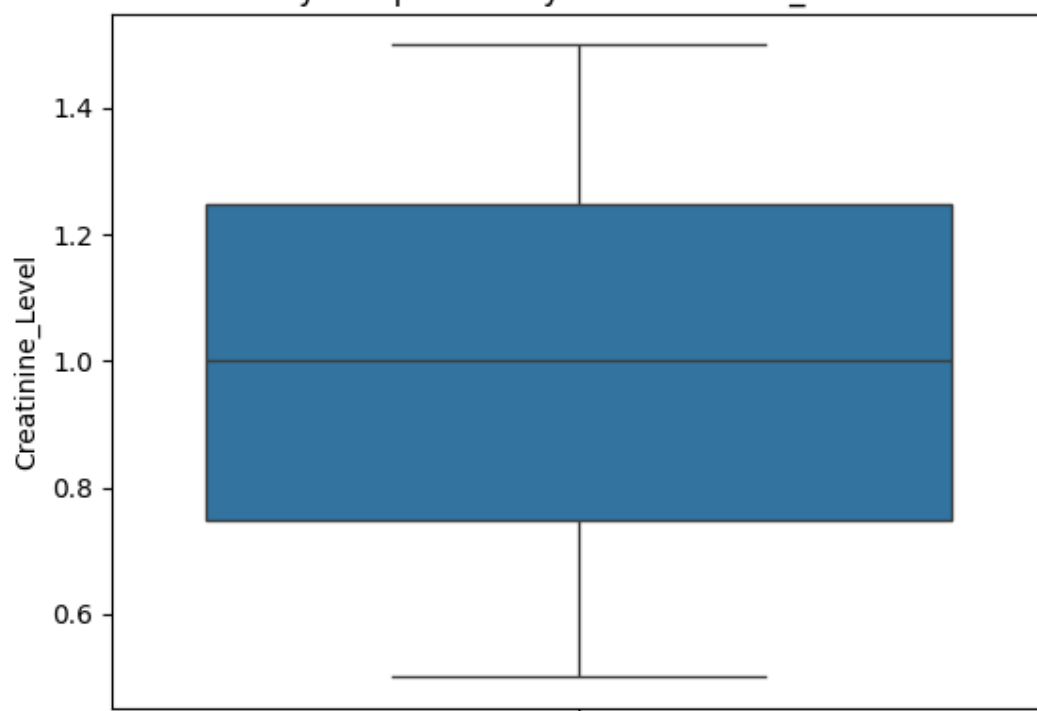
Wykres pudełkowy dla Blood_Pressure_Systolic



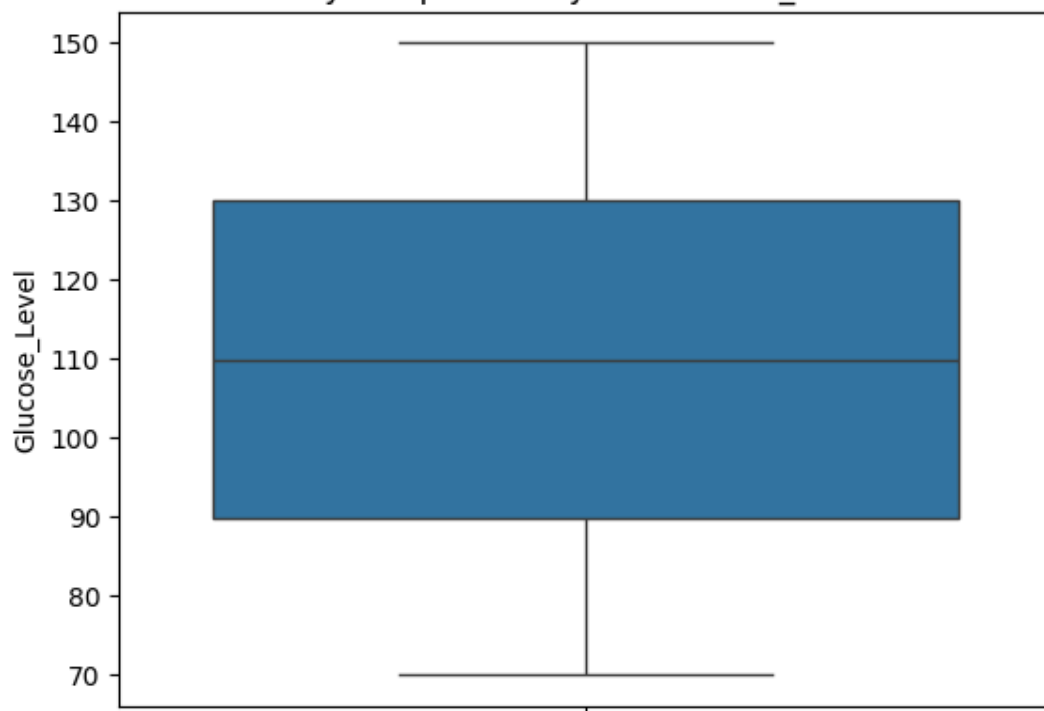
Wykres pudełkowy dla Calcium_Level



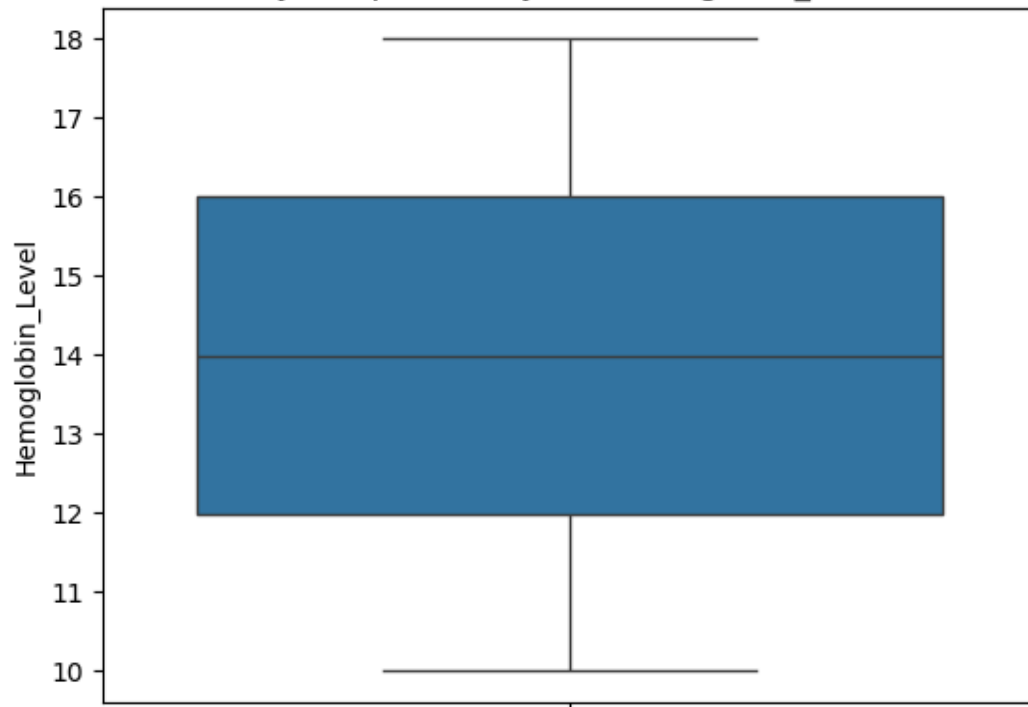
Wykres pudełkowy dla Creatinine_Level



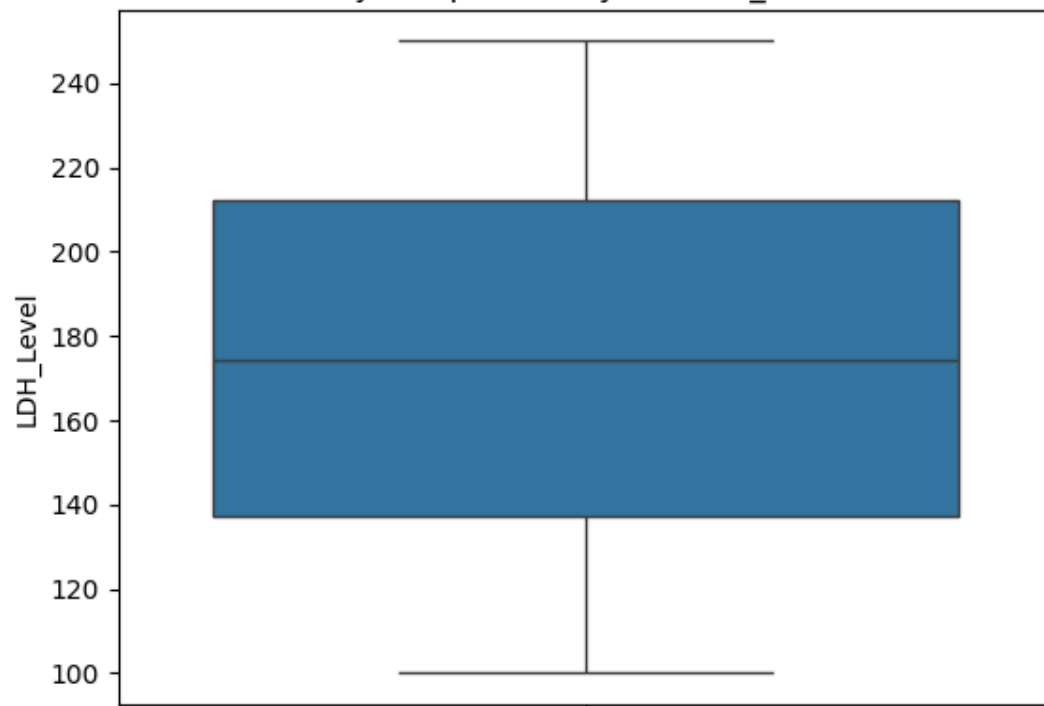
Wykres pudełkowy dla Glucose_Level



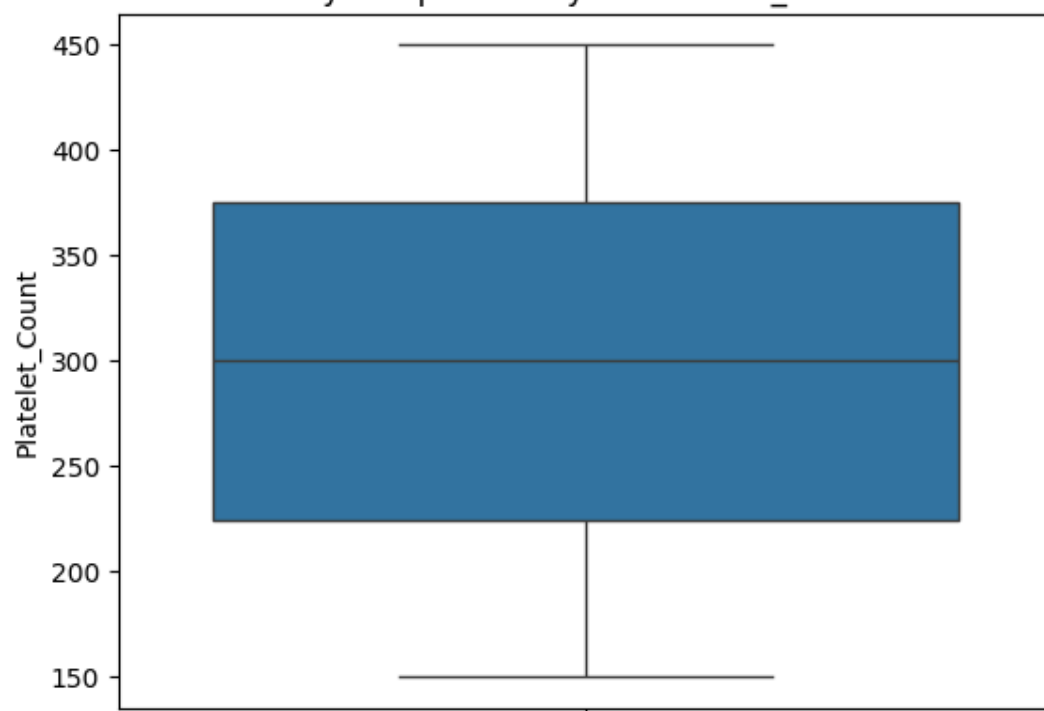
Wykres pudełkowy dla Hemoglobin_Level



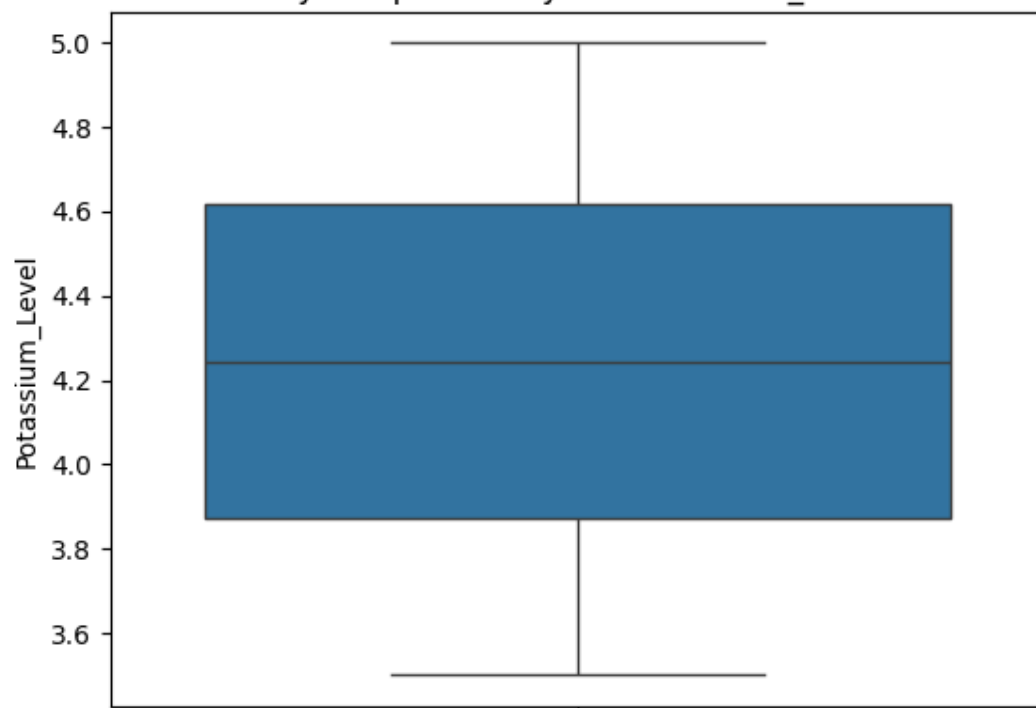
Wykres pudełkowy dla LDH_Level



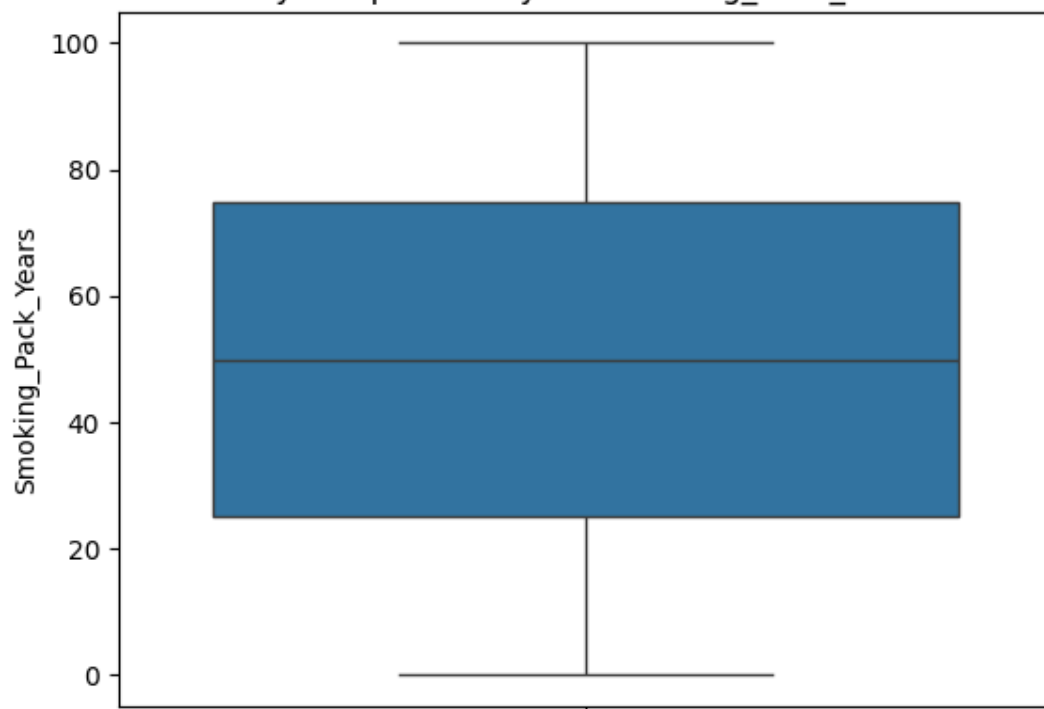
Wykres pudełkowy dla Platelet_Count



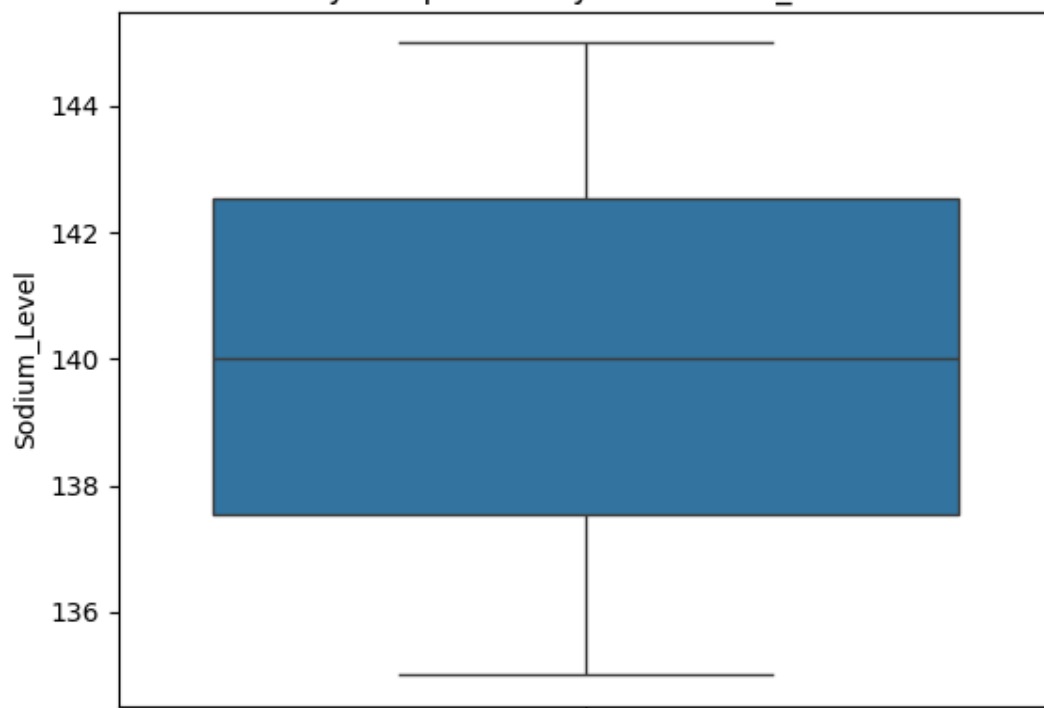
Wykres pudełkowy dla Potassium_Level



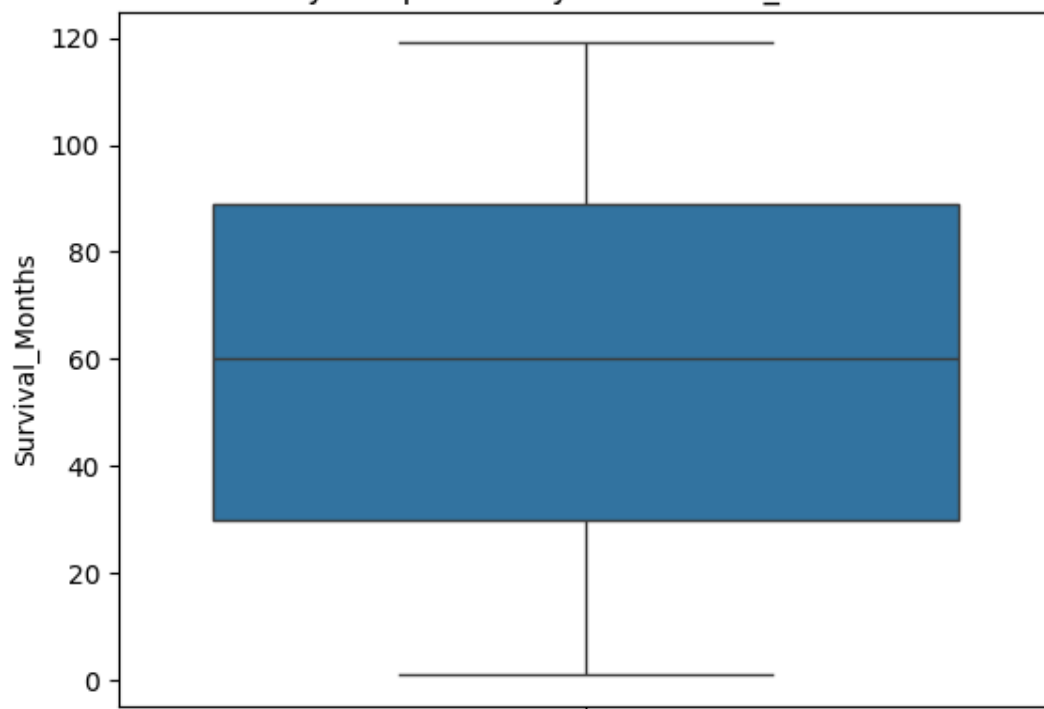
Wykres pudełkowy dla Smoking_Pack_Years



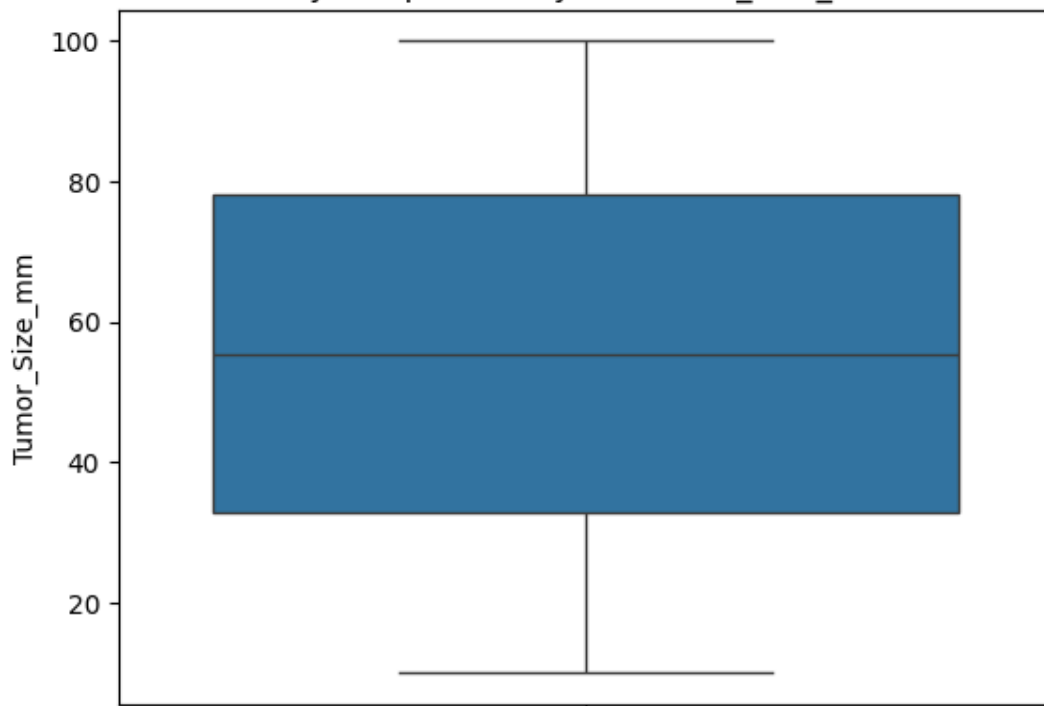
Wykres pudełkowy dla Sodium_Level



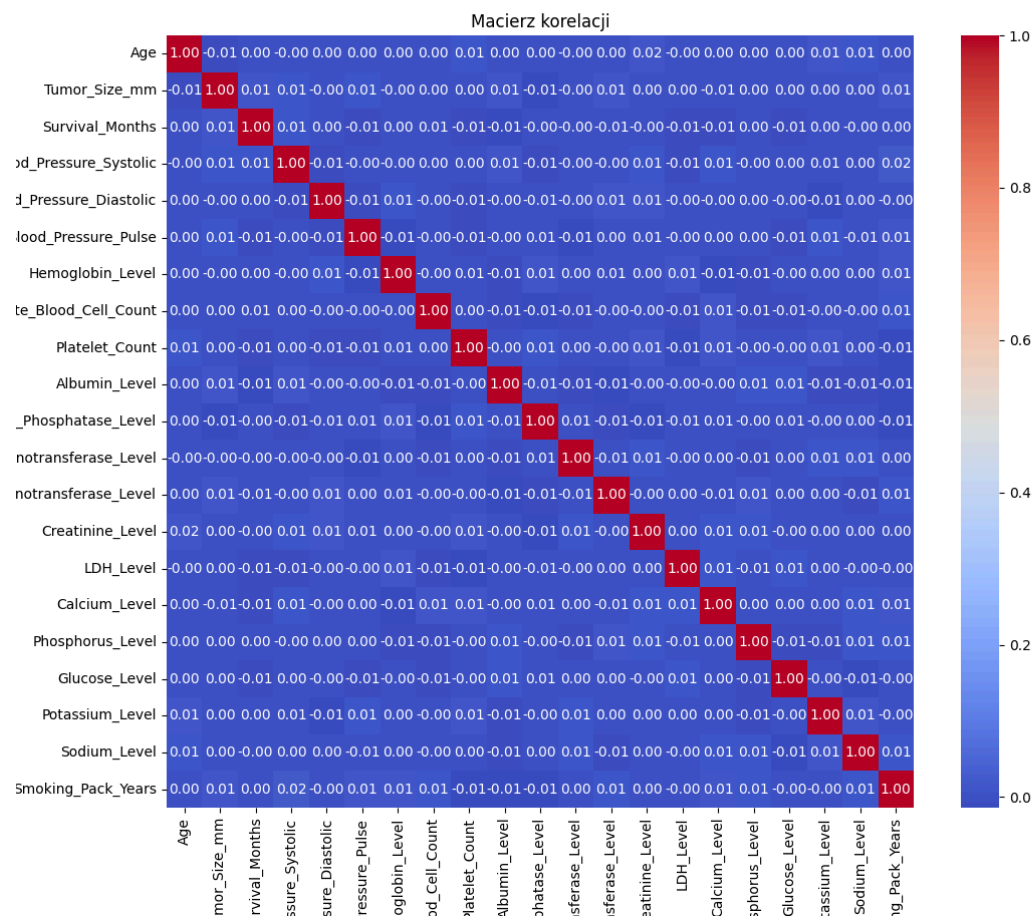
Wykres pudełkowy dla Survival_Months



Wykres pudełkowy dla Tumor_Size_mm



Dalej jest przedstawiona tabela korelacji wartości numerycznych na której można zauważyć, że takich korelacji nie ma.



Na koniec zostało sprawdzone czy są jakieś wartości brakujące, ale ze względu na to, że takowych wartości nie ma to nic oprócz tego nie zostało wyświetlone.

Automatyczna

Niestety analiza zautomatyzowanymi narzędziami okazała się niezbyt możliwa, ze względu na niekompatybilność, jak i za dużą ilość danych przez które nie można było wygenerować wykresów.

Efekty które udało się uzyskać to informacje o rodzaju danych w kolumnach, brakujących wartościach i rekomendacje co do tego co można zrobić z danymi kolumnami. Jedynymi zmianami zaproponowanymi w czasie tej analizy było usunięcie kolumny Patient_ID.

Aktualizacja

Udało się uruchomić Sweetviz po zmianie w pliku graph_numeric.py:

z:

```
warnings.filterwarnings('ignore', category=np.VisibleDeprecationWarning)
```

```
# ...
warnings.filterwarnings('once', category=np.VisibleDeprecationWarning)
```

na:

```
warnings.filterwarnings('ignore', category=DeprecationWarning)
# ...
warnings.filterwarnings('once', category=DeprecationWarning)
```

Ogólnie rzecz biorąc wyniki są podobne do wcześniej uzyskanych.

Cały zbiór danych

Można zobaczyć brak korelacji i równomierne rozłożenie danych jak i czasami dużą różnorodność danych.

Porównanie zbioru danych test i train

Można zobaczyć, że procentowe rozłożenie danych jest dość wyrównane co powinno dobrze wpłynąć na trenowanie modelu.

Podsumowanie analizy

Podsumowując należy usunąć Insurance_Type i Patient_ID, pierwsze ze względu różnicy między krajami, a drugie przez to, że jest to id co nie będzie potrzebne przy tworzeniu modelu. Dodatkowo nie widać pomiędzy kolumnami korelacji, więc trudno będzie wprowadzić jakieś zmiany związane z tym. Jednak można zauważyć wiele wartości które się nie powtarzają przez co można pomyśleć czy nie warto jakoś ich zaokrąglić.

Modele

Bez zaokrąglenia danych

Model	Wynik	Jak został zrobiony	Inne informacje
#1	0.25638457265467657	ExtraTreesClassifier (input_matrix, ExtraTreesClassifier __bootstrap=False, ExtraTreesClassifier __criterion=entropy, ExtraTreesClassifier __max_features=0.4	{'generation': 5, 'mutation_count': 2, 'crossover_count': 0, 'predecessor': ('ExtraTreesClassifie r(input_matrix, ExtraTreesClassifier __bootstrap=False,

		<pre> ExtraTreesClassifier __min_samples_leaf =15, ExtraTreesClassifier __min_samples_spli t=12, ExtraTreesClassifier __n_estimators=100) </pre>	<pre> ExtraTreesClassifier __criterion=entropy, ExtraTreesClassifier __max_features=0.4 , ExtraTreesClassifier __min_samples_leaf =5, ExtraTreesClassifier __min_samples_spli t=12, ExtraTreesClassifier __n_estimators=100),), 'operator_count': 1, 'internal_cv_score': np.float64(0.256384 57265467657)} </pre>
#2	0.25569309463296225	<pre> KNeighborsClassifie r(input_matrix, KNeighborsClassifie r__n_neighbors=82, KNeighborsClassifie r__p=2, KNeighborsClassifie r__weights=uniform) </pre>	<pre> {'generation': 3, 'mutation_count': 1, 'crossover_count': 0, 'predecessor': ('KNeighborsClassifi er(input_matrix, KNeighborsClassifie r__n_neighbors=82, KNeighborsClassifie r__p=2, KNeighborsClassifie r__weights=distance)'), 'operator_count': 1, 'internal_cv_score': np.float64(0.255693 09463296225)} </pre>
#3	0.25509031629399037	<pre> KNeighborsClassifie r(input_matrix, KNeighborsClassifie r__n_neighbors=44, KNeighborsClassifie r__p=1, KNeighborsClassifie r__weights=uniform) </pre>	<pre> {'generation': 5, 'mutation_count': 1, 'crossover_count': 1, 'predecessor': ('BernoulliNB(input_ matrix, BernoulliNB__alpha =100.0, BernoulliNB__fit_pri or=True)'),), 'operator_count': 1, 'internal_cv_score': np.float64(0.255090 31629399037)} </pre>

#4	0.25500392339829914	<p>GradientBoostingClassifier(input_matrix, GradientBoostingClassifier__learning_rate=1.0, GradientBoostingClassifier__max_depth=9, GradientBoostingClassifier__max_features=0.4, GradientBoostingClassifier__min_samples_leaf=19, GradientBoostingClassifier__min_samples_split=3, GradientBoostingClassifier__n_estimators=100, GradientBoostingClassifier__subsample=0.6500000000000001)</p>	<pre>{'generation': 0, 'mutation_count': 0, 'crossover_count': 0, 'predecessor': ('ROOT',), 'operator_count': 1, 'internal_cv_score': np.float64(0.25500392339829914)}</pre>
#5	0.2550036257483571	<p>GradientBoostingClassifier(input_matrix, GradientBoostingClassifier__learning_rate=0.01, GradientBoostingClassifier__max_depth=2, GradientBoostingClassifier__max_features=0.9000000000000001, GradientBoostingClassifier__min_samples_leaf=13, GradientBoostingClassifier__min_samples_split=2, GradientBoostingClassifier__n_estimators=100, GradientBoostingClassifier__subsample=0.7500000000000001)</p>	<pre>{'generation': 1, 'mutation_count': 1, 'crossover_count': 0, 'predecessor': ('BernoulliNB(input_matrix, BernoulliNB__alpha=100.0, BernoulliNB__fit_prior=False)',), 'operator_count': 1, 'internal_cv_score': np.float64(0.2550036257483571)}</pre>

Z zaokrągleniem danych do 2 miejsca po przecinku

Model	Wynik	Jak został zrobiony	Inne informacje
#1	0.2585402279477669	SGDClassifier(input_matrix, SGDClassifier__alpha=0.0, SGDClassifier__eta0=0.01, SGDClassifier__fit_intercept=True, SGDClassifier__l1_ratio=0.25, SGDClassifier__learning_rate=constant, SGDClassifier__loss=modified_huber, SGDClassifier__penalty=elasticnet, SGDClassifier__power_t=10.0)	{'generation': 0, 'mutation_count': 0, 'crossover_count': 0, 'predecessor': ('ROOT',), 'operator_count': 1, 'internal_cv_score': np.float64(0.2585402279477669)}
#2	0.2550028072110163	DecisionTreeClassifier(RBFSampler(input_matrix, RBFSampler__gamma=0.2), DecisionTreeClassifier__criterion=entropy, DecisionTreeClassifier__max_depth=6, DecisionTreeClassifier__min_samples_leaf=11, DecisionTreeClassifier__min_samples_split=5)	{'generation': 4, 'mutation_count': 2, 'crossover_count': 0, 'predecessor': ('DecisionTreeClassifier(RBFSampler(input_matrix, RBFSampler__gamma=0.2), DecisionTreeClassifier__criterion=entropy, DecisionTreeClassifier__max_depth=6, DecisionTreeClassifier__min_samples_leaf=11, DecisionTreeClassifier__min_samples_split=5)'), 'operator_count': 2, 'internal_cv_score': np.float64(0.2550028072110163)}
#3	0.2547440005863704	DecisionTreeClassifier(RBFSampler(input_matrix, RBFSampler__gamma=0.2), DecisionTreeClassifier__criterion=gini, DecisionTreeClassifier__max_depth=6, DecisionTreeClassifier__min_	{'generation': 5, 'mutation_count': 3, 'crossover_count': 0, 'predecessor': ('DecisionTreeClassifier(RBFSampler(input_matrix, RBFSampler__gamma

		<p>samples_leaf=11, DecisionTreeClassifier__min_ samples_split=5)</p>	<p>=0.2), DecisionTreeClassifier __criterion=entropy, DecisionTreeClassifier __max_depth=6, DecisionTreeClassifier __min_samples_leaf=1 1, DecisionTreeClassifier __min_samples_split=5)'), 'operator_count': 2, 'internal_cv_score': np.float64(0.25474400 05863704))}</p>
#4	0.2545712892074735	<p>DecisionTreeClassifier(RBFSa mpler(input_matrix, RBFSampler__gamma=0.2), DecisionTreeClassifier__criteri on=entropy, DecisionTreeClassifier__max_ depth=6, DecisionTreeClassifier__min_ samples_leaf=5, DecisionTreeClassifier__min_ samples_split=11)</p>	<p>{'generation': 4, 'mutation_count': 3, 'crossover_count': 0, 'predecessor': (('DecisionTreeClassifier (RBFSampler(input_ma trix, RBFSampler__gamma =0.2), DecisionTreeClassifier __criterion=entropy, DecisionTreeClassifier __max_depth=6, DecisionTreeClassifier __min_samples_leaf=1 9, DecisionTreeClassifier __min_samples_split=1 1)'),), 'operator_count': 2, 'internal_cv_score': np.float64(0.25457128 92074735))}</p>
#5	0.2543125942015559	<p>KNeighborsClassifier(RobustS caler(input_matrix), KNeighborsClassifier__n_neig hbors=15, KNeighborsClassifier__p=1, KNeighborsClassifier__weight s=uniform)</p>	<p>{'generation': 4, 'mutation_count': 2, 'crossover_count': 0, 'predecessor': (('BernoulliNB(RobustSc aler(input_matrix), BernoulliNB__alpha=1. 0, BernoulliNB__fit_prior= True)'),), 'operator_count': 2, 'internal_cv_score': np.float64(0.25431259 42015559))}</p>

Wybrany model

Wybrałam model #1 z grupy działającej na zaokrąglonych danych, ponieważ daje najlepszy wynik.

Wyniki modelu:

accuracy: 0.2457729468599034

precision: 0.2402012891721489

F1: 0.2190883897895028