Index	df - Ped () ** **Yi
0 1 1 2 2 4 3 13 4 16 5 rows ×	3. 919 2.6999 0 0 0 0 0 0 31.4 2 0 0 0 2.949 1.591 0 7.253 0 0 2 2. 4170 2.1144 0 0 0 0 0 30.8 1 0 0 0 3.315 1.967 0 7.257 0 0 2 3. 4214 2.6272 0 0 0 0 0 30.0 3 0 0 0 2.998 1.722 0 6.770 0 0 2 3. 4214 2.6272 0 0 0 0 31.6 2 0 0 0 3.542 1.739 0 8.127 0 1 2 43 columns Dration Counts = train_df["Class"].value_counts() "Class 1: ", class_counts[1], "percentage: ", class_counts[2]/len(train_df)) "Class 2: ", class_counts[2], "percentage: ", class_counts[2]/len(train_df)) "Class 2: ", class_counts[2], "percentage: ", class_counts[2]/len(train_df)) 1: 564 percentage: 8.666666666666666666666666666666666666
Class 2	1
<pre>nans = fig = fig.su plt.ba</pre>	et variable is binary, with the value 1 indicating that the chemical is bio-degradable and 2 indicating that it is not bio-degradable. The dataset is imbalanced, with 1's representing 66.7% of the data and 2's representing 33.3% of the data. E train_df.isnull().sum(axis = 0) plt.figure(figsize=(19, 5)) ur(nans.index, nans.values) icks(rotation=90) NaNs in columns NaNs in columns
20 - 15 - 10 -	jeg v v v v v v v v v v v v v v v v v v v
correl correl fig = fig.su plt.ba	e a few NaN values in the dataset, but not a lot. We assume that dropping these rows will not have a significant impact on the model, but we will also test the model with imputation such as taking the mean value. Lation_in_data = train_df.corr() Lation_to_class = correlation_in_data["Class"] plt.figure(figsize=(10, 5)) uptitle('Correlation to class vairable', fontsize=16) ur(correlation_to_class.index, correlation_to_class.values) cicks(rotation=90)
0.8 - 0.6 - 0.4 - 0.00.20.40.6 -	1 2 2 2 2 2 2 2 2 2
sns.he AxesSi Index - V2 - V4 - V6 - V10 - V12 - V14 - V16 - V18 - V20 - V24 - V26 - V28 - V28 -	res have a very high direct correlation to the target variable, but quite a lot of features have some correlation. reatmap(correlation_in_data, fmt=".2f") ubplot: > - 1.00 - 0.75 - 0.50 - 0.25 - 0.00 - 0.25
correl thresh	see that most features are not directly correlated to one another, but there are some brighter spots on the heatmap indicating some correlation between features. ated_columns = set()
print(('V38').79882 2), ('V 17 high train_ plt.ba plt.xt plt.ti plt.sh	correlated_columns
plt.xt	dd_without_index_and_class.boxplot(figsize=(10, 10)) itle("Outliers of all columns")
	Outliers of all columns Outliers of all columns
80 - 60 - 20 -	
y plottin contin " ".jc V1 V2 /Ode /e decic Drop Repl	In the distribution of the features, we can see that most features have some outlies. We will test the model with and without outlier removal, we assume that removing the outliers will have a significant impact on the model. In the distribution of the features, we can see that most features have some outlies. We will test the model with and without outlier removal, we assume that removing the outliers will have a significant impact on the model. In the distribution of the features, we can see that most features have some outlies. We will test the model with and without outlier removal, we assume that removing the outliers will have a significant impact on the model. In the distribution of the features, we can see that most features have some outlies. We will test the model with and without outlier removal, we assume that removing the outliers will have a significant impact on the model. In the distribution of the features, we can see that most features have a significant impact on the model. In the distribution of the features have a significant impact on the model. In the distribution of the features have a significant impact on the model. In the distribution of the features have a significant impact on the model. In the distribution of the features have a significant impact on the model. In the distribution of the features have a significant impact on the model. In the distribution of the features have a significant impact on the model. In the distribution of the features have a significant impact on the model. In the distribution of the features have a significant impact on the model. In the distribution of the features have a significant impact on the model. In the distribution of the features have a significant impact on the model. In the distribution of the features have a significant impact on the model. In the distribution of the features have a significant impact on the model. In the distribution of the features have a significant impact on the model. In the distribution of the featur
• Poly train_ withou nan_re withou trans poli_c size = poli_c poli_c test_c test_c test_c	<pre>piping outliers nomial features data = train_df.drop(["Index"], axis=1) tt_nan = train_data_dropna(axis=0) placed = train_data_fillna(without_nan.mean()) tt_outliers = without_nan[(np.abs(stats.zscore(without_nan)) < 3).all(axis=1)] = PolynomialFeatures(degree=3) tata = trans_fit_transform(without_nan.drop(["Class"], axis=1))</pre>
def spread of the spread of th	rest = trans.transform(test_data_without_nan.drop(["class"], axis=1))
//ajori train_ maj_cl best_m //aive wnan_f wnan_t withou mean_f mean_t	rind_ata = train_data.copy() ty Classifier features, train_target = split_data(train_data, "Class") lassifier = Majorityclassifier(train_features, train_target) lassifier = Majorityclassifier(train_features, train_target) las_data = train_data.copy() Bayes Modeling features, wnan_target = split_data(without_nan, "Class") features, wnan_target = split_data(without_nan, "Class") features, wnan_target = split_data(test_data_without_nan, "Class") features, mean_target = split_data(nan_replaced, "Class")
heavy_ outlies outlies outlies poli_f poli_t poli_r nb_clas scores for i,	smoothing_nan_nb = NaiveBayesClassifier(wnan_features, wnan_target, var_smoothing=10e-9) smoothing_nan_nb = NaiveBayesClassifier(wnan_features, wnan_target, var_smoothing=10e-12) ers_features, outliers_target = split_data(without_outliers, "Class") ers_test_features, outliers_test_target = split_data(without_outliers, "Class") ers_nb = NaiveBayesClassifier(outliers_features, outliers_target) features, poli_target = split_data(poli_data, "Class") est_features, poli_test_target = split_data(poli_lest, "Class") ib = NaiveBayesClassifier(poli_features, poli_target) ssifiers = [(without_nan_nb, test_data_without_nan), (mean_nb, test_data_nan_replaced), (slight_smoothing_nan_nb, test_data_without_nan), (heavy_smoothing_nan_nb, test_data_without_nan), (outliers_nb, test_data_sis_sis_sis_sis_sis_sis_sis_sis_sis_si
fig, a score_classifor i ax ax ax ax	cores[i] = c.evaluate(f, t) Lixes = plt.subplots(1, 5, figsize=(15, 5)) Linames = ["Accuracy", "Precision", "Recall", "F1", "AUC"] Infiers = ["Wo Nah", "Mean", "Slight S", "Heavy S", "Outliers", "Poly"] In range(5): Less[i].bar([i for i in range(len(nb_classifiers))], scores[:, i]) Less[i].set.title(score_names[i]) Less[i].set.xitcks(alf for i in range(len(nb_classifiers))]) Less[i].set.xitcks(alf for i in range(len(nb_cla
	ic Regression Modeling
mean_fmean_l L1_per balance outlie outlie	Teatures, wnan_target = split_data(without_nan, "Class") tr_nan_lr = LogisticRegressionClassifier(wnan_features, wnan_target, solver='lbfgs', max_iter=1000) Teatures, mean_target = split_data(nan_replaced, "Class") tr = LogisticRegressionClassifier(mean_features, mean_target, solver='lbfgs', max_iter=1000) malty_lr = LogisticRegressionClassifier(mean_features, mean_target, solver='lbfgs', max_iter=1000, penalty='l1') malty_lr = LogisticRegressionClassifier(mean_features, mean_target, solver='lbfgs', max_iter=1000, class_weight='balanced')
 ز_oli_ن	Features, poli_target = split_data(poli_data, "Class") Lr = LogisticRegressionClassifier(poli_features, poli_target, solver='lbfgs', max_iter=1000, tol=1e-2)
poli_l c:\User STOP: 1 Increase htt Please htt n_ite fig, a score_ classi for i ax ax ax	features, poli_target = split_data(poli_data, "Class") for = LogisticRegressionClassifier(poli_features, poli_target, solver='lbfgs', max_iter=1000, tol=1e-2) rs\Aleksander\anaconda3\envs\inteligent_systems\lib\site-packages\sklearn\linear_model_logistic.py:444: ConvergenceWarning: lbfgs failed to converge (status=1): rotal NO. of ITERATIONS REACHED LIMIT. see the number of iterations (max_iter) or scale the data as shown in: tps://scikit-learn.org/stable/modules/preprocessing.html also refer to the documentation for alternative solver options: tps://scikit-learn.org/stable/modules/linear_model.html#logistic-regression tps://scikit-learn.org/stable/modules/linear_model.html#logistic-regression tpr_i = _check_optimize_result(
fig, a score classifor i ax	Teatures, poli_target = split_data(poli_data, "class") r = logisticRegressionClassifier(poli_features, poli_target, solver='lbfgs', max_iter=1000, tol=1e-2) rr = logisticRegressionClassifier(poli_features, poli_target, solver='lbfgs', max_iter=1000, tol=1e-2) roTAL MO. of ITERATIONS REACHED LIMIT. see the number of iterations (max_iter) or scale the data as shown in: tps://scikit-learn.org/stable/modules/preprocessing.html also refer to the documentation for alternative solver options: tps://scikit-learn.org/stable/modules/linear_model.html#logistic-regression er_i = _check_optimize_result(sesifiers = [(without_nan_lr, test_data_without_nan), (mean_lr, test_data_nan_replaced), (Li_penalty_lr, test_data_nan_replaced), (balanced_lr, test_data_nan_replaced), (outliers_lr, test_data_without_outliers), takes = plt.subplots(1, 5, figsize=(15, 5)) names = ["Accuracy", "Precision", "Recall", "F1", "AUC"] fiers = ["www.Nah", "Nean", "L1", "Balanced", "Outliers", "Poly"] in range(5): tes[i].set_ticle(score_names[i]) tes[i].set_ticle(score_names[i]) tes[i].set_ticle(score_names[i]) tes[i].set_ticls(score_names[i]) tes[i].set_ticls(slore_names[i])
poli_l c:\User STOP: Tolease htt n_ite lease htt n_ite lr_cla fig, a score classi for i ax ax ax ax 0.8 - 0.7 - 0.6 - 0.7 - 0.6 - 0.7 - 0.7 - 0.7 - 0.7 - 0.8 - 0.9 - 0.9 - 0.1 - 0.1 - 0.1 - 0.2 - 0.1 - 0.3 - 0.4 - 0.5 - 0.6 - 0.7 - 0.7 - 0.7 - 0.8 - 0.9 - 0.9 - 0.9 - 0.1 - 0.1 - 0.1 - 0.2 - 0.1 - 0.3 - 0.3 - 0.4 - 0.5 - 0.6 - 0.7 - 0.7 - 0.8 - 0.9 - 0.9 - 0.9 - 0.1 - 0.1 - 0.1 - 0.2 - 0.3 - 0.3 - 0.4 - 0.5 - 0.6 - 0.7 - 0.7 - 0.7 - 0.8 - 0.9 -	re unjustichengressientlessifier (polifications, molitarget, solver='lbfgs', max_iter=1890, tol=12:2) **SANDERSORD'S MOLITAGE MOLIDATION (Color LINT) **SANDERSORD'S MOLITAGE
poli_I C:\User GTOP: I Increase Increas	Particle (Control of Special Action), and import otherwises and production there are subject to the control of Special Action (Control of Special Action), and import otherwises and production there are subject to the control of Special Action (Control of Special Action) (Control of Special Action of
poli_1 C:\User GTOP: Increase	The properties of the properti
poli_i croreas fig, a score classi for i ax ax ax 0.8 0.7 0.6 0.7 0.6 0.7 0.7 0.7 0.7	The property of the property o
poli_i c: User TOP: I chose htt flease ht flease htt flease ht fl	The process of the control of the co
poli_i C:\User STOP: I Coreas fig, a ax ax ax ax ax ax ax ax ax	
poli_l c:\User STOP: I c:\User STOP: I Increas	A THE MACHINE AND
Please OR OR OR OR OR OR OR OR OR O	The control of the co