

## Imports

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
from scipy import stats
from sklearn.preprocessing import RobustScaler
from sklearn.preprocessing import PolynomialFeatures

from classifiers.random_classifier import RandomClassifier
from classifiers.majority_classifier import MajorityClassifier
from classifiers.naive_bayes_classifier import NaiveBayesClassifier
from classifiers.logistic_regression_classifier import LogisticRegressionClassifier
from classifiers.decision_tree_classifier import DecisionTreeClassifier

train_df = pd.read_csv("../seminar_2/train.csv")

train_df.head()

Index      V1      V2      V3      V4      V5      V6      V7      V8      V9      ...      V33      V34      V35      V36      V37      V38      V39      V40      V41      Class
0      1.3919      2.6909      0      0      0      0      0      31.4      2      ...      0      0      0      2.949      1.591      0      7.257      0      0      2
1      5.4226      3.2544      0      0      0      0      0      29.4      2      ...      0      0      0      3.391      2.405      0      6.003      0      0      2
2      6.4226      3.4296      0      0      0      0      0      29.6      2      ...      0      0      0      3.351      2.596      0      7.904      0      0      2
3      7.5000      5.0476      1      0      0      0      0      11.1      0      ...      0      0      0      1.472      4.583      0      9.303      0      0      2
4      8.4525      3.8301      0      0      0      0      0      31.6      3      ...      0      0      0      3.379      2.143      0      7.950      0      0      2

5 rows x 43 columns

test_df = pd.read_csv("../seminar_2/test.csv")

test_df.head()

Index      V1      V2      V3      V4      V5      V6      V7      V8      V9      ...      V33      V34      V35      V36      V37      V38      V39      V40      V41      Class
0      1.3919      2.6909      0      0      0      0      0      31.4      2      ...      0      0      0      2.949      1.591      0      7.257      0      0      2
1      5.4226      3.2544      0      0      0      0      0      29.4      2      ...      0      0      0      3.391      2.405      0      6.003      0      0      2
2      6.4226      3.4296      0      0      0      0      0      29.6      2      ...      0      0      0      3.351      2.596      0      7.904      0      0      2
3      7.5000      5.0476      1      0      0      0      0      11.1      0      ...      0      0      0      1.472      4.583      0      9.303      0      0      2
4      8.4525      3.8301      0      0      0      0      0      31.6      3      ...      0      0      0      3.379      2.143      0      7.950      0      0      2

5 rows x 43 columns
```

## Exploration

```
class_counts = train_df["Class"].value_counts()

print("Class 1: ", class_counts[1], "percentage: ", class_counts[1]/len(train_df))
print("Class 2: ", class_counts[2], "percentage: ", class_counts[2]/len(train_df))

Class 1: 564 percentage: 0.6666666666666666
Class 2: 282 percentage: 0.3333333333333333

p = plt.pie(class_counts, labels=class_counts.index, autopct='%1.1f%%')

1
66.7%
33.3%
2
```

The target 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.

```
nans = train_df.isnull().sum(axis=0)

fig = plt.figure(figsize=(10, 5))
fig.suptitle('Nans in columns', fontsize=16)
plt.bar(train_df.index, nans.values)
plt.xticks(rotation=90)
plt.show()

NaNs in columns

Index      V1      V2      V3      V4      V5      V6      V7      V8      V9      ...      V33      V34      V35      V36      V37      V38      V39      V40      V41      Class
0      0      0      0      0      0      0      0      0      0      ...      0      0      0      0      0      0      0      0      0      0
1      0      0      0      0      0      0      0      0      0      ...      0      0      0      0      0      0      0      0      0      0
2      0      0      0      0      0      0      0      0      0      ...      0      0      0      0      0      0      0      0      0      0
3      0      0      0      0      0      0      0      0      0      ...      0      0      0      0      0      0      0      0      0      0
4      0      0      0      0      0      0      0      0      0      ...      0      0      0      0      0      0      0      0      0      0

There are 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.
```

```
correlation_in_data = train_df.corr()

correlation_to_class = correlation_in_data["Class"]

fig = plt.figure(figsize=(10, 5))
fig.suptitle('Correlation to class variable', fontsize=16)
plt.bar(correlation_in_data.index, correlation_to_class.values)
plt.xticks(rotation=90)
plt.show()

Correlation to class variable

Index      V1      V2      V3      V4      V5      V6      V7      V8      V9      ...      V33      V34      V35      V36      V37      V38      V39      V40      V41      Class
0      -0.5      0.1      0.1      0.1      0.1      0.1      0.1      0.1      0.1      ...      0.1      0.1      0.1      0.1      0.1      0.1      0.1      0.1      0.1      0.1
1      -0.5      0.1      0.1      0.1      0.1      0.1      0.1      0.1      0.1      ...      0.1      0.1      0.1      0.1      0.1      0.1      0.1      0.1      0.1      0.1
2      -0.5      0.1      0.1      0.1      0.1      0.1      0.1      0.1      0.1      ...      0.1      0.1      0.1      0.1      0.1      0.1      0.1      0.1      0.1      0.1
3      -0.5      0.1      0.1      0.1      0.1      0.1      0.1      0.1      0.1      ...      0.1      0.1      0.1      0.1      0.1      0.1      0.1      0.1      0.1      0.1
4      -0.5      0.1      0.1      0.1      0.1      0.1      0.1      0.1      0.1      ...      0.1      0.1      0.1      0.1      0.1      0.1      0.1      0.1      0.1      0.1

No features have a very high direct correlation to the target variable, but quite a lot of features have some correlation.
```

```
sns.heatmap(correlation_in_data, fctg="2*")

<AxesSubplot: >

Index      V1      V2      V3      V4      V5      V6      V7      V8      V9      ...      V33      V34      V35      V36      V37      V38      V39      V40      V41      Class
0      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0      ...      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0
1      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0      ...      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0
2      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0      ...      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0
3      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0      ...      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0
4      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0      ...      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0      0.0

We can 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.
```

```
correlated_columns = set()
threshold = 0.75

for i in range(len(correlation_in_data.columns)):
    for j in range(i+1, len(correlation_in_data.columns)):
        if abs(correlation_in_data.iloc[i, j]) > threshold:
            colnames = correlation_in_data.columns[i]
            colname2 = correlation_in_data.columns[j]
            correlated_columns.add((colnames, colname2, correlation_in_data.iloc[i, j]))

print(correlated_columns)
print(len(correlated_columns), "highly correlated features")

(('V38', 'V11', 0.836978414216805), ('V27', 'V1', 0.92156062534691), ('V39', 'V36', 0.9165966183518999), ('V38', 'V18', 0.7577897873018442), ('V37', 'V17', 0.8498222432821918), ('V39', 'V13', 0.813370360627613), ('V38', 'V34', 0.788841728566266), ('V22', 'V18', -0.8008371258956488), ('V33', 'V7', 0.788841728566266), ('V16', 'V18', 0.843963443813359), ('V15', 'V1', 0.8072060608897), ('V11', 'V5', 0.850815645318726), ('V29', 'V28', 0.844062668784212), ('V34', 'V5', 0.732712220086253), ('V34', 'V11', 0.7832615815925773), ('V2', 'V39', 0.90134043848) features
```

```
train_df_without_index_and_class = train_df.drop(["Index", "Class"], axis=1)
plt.bar(train_df_without_index_and_class.nunique().index, train_df_without_index_and_class.nunique().values)
plt.xticks(rotation=90)

Number of unique values in columns

Index      V1      V2      V3      V4      V5      V6      V7      V8      V9      ...      V33      V34      V35      V36      V37      V38      V39      V40      V41      Class
0      1      2      3      4      5      6      7      8      9      ...      10      11      12      13      14      15      16      17      18      19
1      1      2      3      4      5      6      7      8      9      ...      10      11      12      13      14      15      16      17      18      19
2      1      2      3      4      5      6      7      8      9      ...      10      11      12      13      14      15      16      17      18      19
3      1      2      3      4      5      6      7      8      9      ...      10      11      12      13      14      15      16      17      18      19
4      1      2      3      4      5      6      7      8      9      ...      10      11      12      13      14      15      16      17      18      19

Text(0.5, 1.0, 'Outliers of all columns')
```

```
Text(0.5, 1.0, 'Outliers of all columns')

By plotting the distribution of the features, we can see that most features have some outliers. We will test the model with and without outlier removal, we assume that removing the outliers will have a significant impact on the model.
```

```
continuous_columns = [i for i in train_df_without_index_and_class.nunique().index where (train_df_without_index_and_class.nunique().values >= 100) if i is not None]
nb_continuous_columns

"V1 V2 V8 V12 V13 V14 V15 V17 V18 V22 V27 V28 V30 V31 V36 V37 V39"

Modeling

We decided that we will test models with differently preprocessed data to see which preprocessing method works best. We will test the following preprocessing methods:
```

- Dropping NaN values
- Replacing NaN values with the mean value
- Dropping outliers
- Polynomial features

```
train_data = train_df.drop(["Index"], axis=1)
without_nan = train_data.dropna(axis=0)
nan_replaced = train_data.fillna(without_nan.mean())
without_outliers = without_nan[np.abs(stats.zscore(without_nan)) < 3].all(axis=1)

trans = PolynomialFeatures(degree=3)
poli_data = trans.fit_transform(without_nan.drop(["Class"], axis=1))
size = poli_data.shape[1]
scores = np.empty((poli_data, without_nan["Class"].values.reshape(-1, 1)))
poli_data = pd.DataFrame(poli_data, columns=["poly_{}".format(i) for i in range(size)] + ["Class"])

test_data = test_df.drop(["Index"], axis=1)
test_data_without_nan = test_data.dropna(axis=0)
test_data_nan_replaced = test_data.fillna(test_data_without_nan.mean())
test_data_without_outliers = test_data_without_nan
```

```
poli_test = trans.transform(test_data_without_nan.drop(["Class"], axis=1))
size = poli_test.shape[1]
poli_test = np.hstack((poli_test, test_data_without_nan["Class"].values.reshape(-1, 1)))
poli_test = pd.DataFrame(poli_test, columns=["poly_{}".format(i) for i in range(size)] + ["Class"])

def split_data(data, target_column):
    return data.drop([target_column], axis=1), data[target_column]

We have decided to test the following models:
```

- Random classifier (as a baseline)
- Majority classifier (as a baseline)
- Naive Bayes classifier (because it is fast and simple)
- Logistic regression (because it is good for binary classification)
- Decision tree (because it is good for high dimensional data)

### Random Classifier

```
train_features, train_target = split_data(train_data, "Class")
best_rnd = RandomClassifier(train_features, train_target)
best_rnd_data = train_data.copy()
```

### Majority Classifier

```
train_features, train_target = split_data(train_data, "Class")
maj_classifier = MajorityClassifier(train_features, train_target)
best_maj_data = train_data.copy()
```

## Naive Bayes Modeling

```
wnan_features, wnan_target = split_data(without_nan, "Class")
wnan_test_features, wnan_test_target = split_data(test_data_without_nan, "Class")
without_nan_nb = NaiveBayesClassifier(wnan_features, wnan_target)

mean_features, mean_target = split_data(nan_replaced, "Class")
mean_test_features, mean_test_target = split_data(test_data_nan_replaced, "Class")
mean_nb = NaiveBayesClassifier(mean_features, mean_target)

slight_smoothing_nan_nb = NaiveBayesClassifier(wnan_features, wnan_target, var_smoothing=10e-9)

heavy_smoothing_nan_nb = NaiveBayesClassifier(wnan_features, wnan_target, var_smoothing=10e-12)

outliers_features, outliers_target = split_data(without_outliers, "Class")
outliers_test_features, outliers_test_target = split_data(test_data_without_outliers, "Class")
outliers_nb = NaiveBayesClassifier(outliers_features, outliers_target)

poli_features, poli_target = split_data(poli_data, "Class")
poli_test_features, poli_test_target = split_data(poli_test, "Class")
poli_nb = NaiveBayesClassifier(poli_features, poli_target)

nb_classifiers = [(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_outliers)]

scores = np.empty((len(nb_classifiers), 5))
for i, (c, test_data) in enumerate(nb_classifiers):
    f, t = split_data(test_data, "Class")
    scores[i] = c.evaluate(f, t)

fig, axes = plt.subplots(1, 5, figsize=(15, 5))

score_names = ["Accuracy", "Precision", "Recall", "F1", "AUC"]
classifiers = ["No NaN", "Mean", "L1", "Balanced", "Outliers", "Poly"]
for i in range(5):
    axes[i].bar([i for i in range(len(nb_classifiers))], scores[:, i])
    axes[i].set_title(score_names[i])
    axes[i].set_xticks([i for i in range(len(nb_classifiers))])
    axes[i].set_xticklabels(classifiers)
    axes[i].axis.set_tick_params(rotation=90)
```

Accuracy Precision Recall F1 AUC

Wo NaN Mean L1 Balanced Outliers Poly

### Logistic Regression Modeling

```
wnan_features, wnan_target = split_data(without_nan, "Class")
mean_lr = LogisticRegressionClassifier(wnan_features, wnan_target, solver='lbfgs', max_iter=1000)

mean_features, mean_target = split_data(nan_replaced, "Class")
mean_lr = LogisticRegressionClassifier(mean_features, mean_target, solver='lbfgs', max_iter=1000)

L1_penalty_lr = LogisticRegressionClassifier(mean_features, mean_target, solver='l1linear', max_iter=1000, penalty='l1')

balanced_lr = LogisticRegressionClassifier(mean_features, mean_target, solver='lbfgs', max_iter=1000, class_weight='balanced')

outliers_features, outliers_target = split_data(without_outliers, "Class")
outliers_lr = LogisticRegressionClassifier(outliers_features, outliers_target, solver='lbfgs', max_iter=2000)

poli_features, poli_target = split_data(poli_data, "Class")
poli_lr = LogisticRegressionClassifier(poli_features, poli_target, solver='lbfgs', max_iter=1000, tol=1e-2)

lr_classifiers = [(without_nan_lr, test_data_without_nan), (mean_lr, test_data_nan_replaced), (L1_penalty_lr, test_data_nan_replaced), (balanced_lr, test_data_nan_replaced), (outliers_lr, test_data_without_outliers)]

fig, axes = plt.subplots(1, 5, figsize=(15, 5))

score_names = ["Accuracy", "Precision", "Recall", "F1", "AUC"]
classifiers = ["No NaN", "Mean", "L1", "Balanced", "Outliers", "Poly"]
for i in range(5):
    axes[i].bar([i for i in range(len(lr_classifiers))], scores[:, i])
    axes[i].set_title(score_names[i])
    axes[i].set_xticks([i for i in range(len(lr_classifiers))])
    axes[i].set_xticklabels(classifiers)
    axes[i].axis.set_tick_params(rotation=90)
```

Accuracy Precision Recall F1 AUC

Wo NaN Mean L1 Balanced Outliers Poly

### Decision Tree Modeling

```
wnan_features, wnan_target = split_data(without_nan, "Class")
without_nan_dt = DecisionTreeClassifier(wnan_features, wnan_target, random_state=42)

mean_features, mean_target = split_data(nan_replaced, "Class")
mean_dt = DecisionTreeClassifier(mean_features, mean_target, random_state=42)

limited_df = DecisionTreeClassifier(mean_features, mean_target, random_state=42, max_depth=10)

cc_df = DecisionTreeClassifier(mean_features, mean_target, random_state=42, criterion='entropy', ccp_alpha=0.01)

outliers_features, outliers_target = split_data(without_outliers, "Class")
outliers_dt = DecisionTreeClassifier(outliers_features, outliers_target, random_state=42)

poli_features, poli_target = split_data(poli_data, "Class")
poli_dt = DecisionTreeClassifier(poli_features, poli_target, random_state=42)

dt_classifiers = [(without_nan_dt, test_data_without_nan), (mean_dt, test_data_nan_replaced), (limited_df, test_data_nan_replaced), (cc_df, test_data_nan_replaced), (outliers_dt, test_data_without_outliers)]

fig, axes = plt.subplots(1, 5, figsize=(15, 5))

score_names = ["Accuracy", "Precision", "Recall", "F1", "AUC"]
classifiers = ["No NaN", "Mean", "L1", "Balanced", "Outliers", "Poly"]
for i in range(5):
    axes[i].bar([i for i in range(len(dt_classifiers))], scores[:, i])
    axes[i].set_title(score_names[i])
    axes[i].set_xticks([i for i in range(len(dt_classifiers))])
    axes[i].set_xticklabels(classifiers)
    axes[i].axis.set_tick_params(rotation=90)
```

Accuracy Precision Recall F1 AUC

Wo NaN Mean L1 Balanced Outliers Poly

### Best Models

```
best_models = [(best_rnd, best_rnd_data), (maj_classifier, best_maj_data), (outliers_nb, test_data_without_outliers), (outliers_lr, test_data_without_outliers), (outliers_dt, test_data_without_outliers)]

len(best_models)

5

As we can see from testing of the models above removing the outliers had the most significant impact on the model. We will further test these models using folding and multiple runs to see if the results are consistent.

### Evaluation



```
repetitions = 10
folds = 5
evaluations = 5
scores = np.empty((model, len(best_models), repetitions, evaluations))
for i, (classifier, model) in enumerate(best_models):
    scores[i] = classifier.test(model_data, "Class", folds=folds, repetitions=repetitions)

score_names = ["F1 score", "Precision", "Recall", "Area under ROC curve", "Accuracy"]

fig, ax = plt.subplots(2, 3, figsize=(15, 10))
for i in range(len(score_names)):
    ax[i, 0].set_title(score_names[i])
    ax[i, 0].bar([i for i in range(len(scores))], scores.mean(axis=(2, 3)), color=sns.color_palette("Set2", 10))
    ax[i, 0].axis.set_ylim(0, 1)
    ax[i, 0].axis.set_xlim(0, 1)

ax[0, 0].set_ylabel("F1 score")
ax[0, 0].set_xlabel("Model")
ax[0, 0].axis.set_xlim(0, 1)
ax[0, 0].axis.set_ylim(0, 1)

ax[0, 1].set_ylabel("Precision")
ax[0, 1].set_xlabel("Model")
ax[0, 1].axis.set_xlim(0, 1)
ax[0, 1].axis.set_ylim(0, 1)

ax[0, 2].set_ylabel("Recall")
ax[0, 2].set_xlabel("Model")
ax[0, 2].axis.set_xlim(0, 1)
ax[0, 2].axis.set_ylim(0, 1)

ax[0, 3].set_ylabel("Area under ROC curve")
ax[0, 3].set_xlabel("Model")
ax[0, 3].axis.set_xlim(0, 1)
ax[0, 3].axis.set_ylim(0, 1)

ax[0, 4].set_ylabel("Accuracy")
ax[0, 4].set_xlabel("Model")
ax[0, 4].axis.set_xlim(0, 1)
ax[0, 4].axis.set_ylim(0, 1)

ax[1, 0].set_ylabel("F1 score")
ax[1, 0].set_xlabel("Model")
ax[1, 0].axis.set_xlim(0, 1)
ax[1, 0].axis.set_ylim(0, 1)

ax[1, 1].set_ylabel("Precision")
ax[1, 1].set_xlabel("Model")
ax[1, 1].axis.set_xlim(0, 1)
ax[1, 1].axis.set_ylim(0, 1)

ax[1, 2].set_ylabel("Recall")
ax[1, 2].set_xlabel("Model")
ax[1, 2].axis.set_xlim(0, 1)
ax[1, 2].axis.set_ylim(0, 1)

ax[1, 3].set_ylabel("Area under ROC curve")
ax[1, 3].set_xlabel("Model")
ax[1, 3].axis.set_xlim(0, 1)
ax[1, 3].axis.set_ylim(0, 1)

ax[1, 4].set_ylabel("Accuracy")
ax[1, 4].set_xlabel("Model")
ax[1, 4].axis.set_xlim(0, 1)
ax[1, 4].axis.set_ylim(0, 1)

We chose logistic regression because it performs best in most metrics (but not by a large margin), but we also decided to test Naive Bayes because for our problem Precision is very important, we do not want to classify non-biodegradable chemicals as biodegradable. Naive Bayes performed better in Precision metric, but worse in Recall metric.
```


```

```
lr_classifier, lr_train = best_models[3]
nb_classifier, nb_train = best_models[2]

train_features, train_target = split_data(lr_train, "Class")

lr_classifier.fit(train_features, train_target)
nb_classifier.evaluate(train_features, train_target)

(0.907037037037037,
 0.9007037037037037,
 0.9007037037037037,
 0.8561461461461461,
 0.87508980144819)

nb_classifiers.fit(train_features, train_target)
nb_classifiers.evaluate(train_features, train_target)

(0.7942456149350678,
 0.978494823695914,
 0.8749187491874919,
 0.8235235235235235,
 0.77564362206987)

Running the models on test data we can see that they performed as expected. Logistic regression performed better in all metrics, but precision where Naive Bayes performed better achieving a score of 0.97.
```