





```

#python-input-38-a246f41ad3640d:8: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
data_for_predictions[col] = ss.transform(data_for_predictions[col].to_numpy().reshape(-1,1))
#python-input-38-a246f41ad3640d:8: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
data_for_predictions[col] = mm.fit_transform(data_for_predictions[col].to_numpy().reshape(-1,1)).reshape(1,-1)

#python-input-38-a246f41ad3640d:8: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
data_for_predictions[col] = mm.fit_transform(data_for_predictions[col].to_numpy().reshape(-1,1)).reshape(1,-1)

#python-input-38-a246f41ad3640d:8: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
data_for_predictions[col] = mm.fit_transform(data_for_predictions[col].to_numpy().reshape(-1,1)).reshape(1,-1)

#python-input-38-a246f41ad3640d:8: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
data_for_predictions[col] = mm.fit_transform(data_for_predictions[col].to_numpy().reshape(-1,1)).reshape(1,-1)

In [39]:
data_for_predictions.head()

Out[39]:
   feature_1  feature_2  feature_3  feature_133  feature_345  feature_954  feature_1044  feature_1048  feature_1110  feature_1137  ...
0  0.154930  0.245826  0.026642             0             0             0             0             0             0             0  ...
1  0.029734  0.128326  0.069143             0             0             0             0             0             0             0  ...
2  0.153935  0.031777  0.025227             0             0             0             0             0             0             0  ...
3  0.061033  0.081033  0.016642             0             0             0             0             0             0             0  ...
4  0.071714  0.364632  0.324983             0             0             0             0             0             0             0  ...

5 rows x 27 columns

In [40]:
data_for_predictions.describe()

Out[40]:
   count      feature_1  feature_2  feature_3  feature_133  feature_345  feature_954  feature_1044  feature_1048  feature_1110  feature_1137  ...
0      5  0.1549300000  0.2458260000  0.0266420000  0.00000000  0.00000000  0.00000000  0.00000000  0.00000000  0.00000000  0.00000000  ...
1      5  0.0297340000  0.1283260000  0.0691430000  0.00000000  0.00000000  0.00000000  0.00000000  0.00000000  0.00000000  0.00000000  ...
2      5  0.1539350000  0.0317770000  0.0252270000  0.00000000  0.00000000  0.00000000  0.00000000  0.00000000  0.00000000  0.00000000  ...
3      5  0.0610330000  0.0810330000  0.0166420000  0.00000000  0.00000000  0.00000000  0.00000000  0.00000000  0.00000000  0.00000000  ...
4      5  0.0717140000  0.3646320000  0.3249830000  0.00000000  0.00000000  0.00000000  0.00000000  0.00000000  0.00000000  0.00000000  ...

8 rows x 27 columns

# ## Method : Following the conventional method

Split the data

In [41]:
X, y = data_for_predictions.iloc[:,1:], data_for_predictions.iloc[:,~1]
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=42)

Libraries for the models

In [42]:
from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier, ExtraTreeClassifier
from sklearn.tree import DecisionTreeClassifier

#model selection
from sklearn.model_selection import GridSearchCV, cross_val_score

#metrics
from sklearn.metrics import confusion_matrix, roc_curve

Algorithms

In [43]:
parameters = {'penalty':['l1','l2'],
              'C': [0.001, 0.01, 0.1, 1]}

#Logistic Regression
GD_LR = GridSearchCV(LogisticRegression(solver='liblinear'),
                    param_grid=parameters,
                    scoring='f1')

#Decision Trees
parameters = {'max_leaf_nodes': [i for i in range(2,10)],
              'max_depth': [i for i in range(2,10)]}

GD_DT = GridSearchCV(DecisionTreeClassifier(),
                    param_grid=parameters,
                    scoring='f1')

#Extra Trees
GD_ET = GridSearchCV(ExtraTreeClassifier(class_weight='balanced'),
                    param_grid=parameters,
                    scoring='f1')

#Random Forest
GD_RF = GridSearchCV(RandomForestClassifier(class_weight='balanced'),
                    param_grid=parameters,
                    scoring='f1')

#Gradient Boosting
GD_GB = GridSearchCV(GradientBoostingClassifier(),
                    param_grid=parameters,
                    scoring='f1')

#K-NN
parameters = {'n_neighbors': [i for i in range(5,35,5)]}
GD_KNN = GridSearchCV(KNeighborsClassifier(),
                    param_grid=parameters,
                    scoring='f1')

models = {'Logistic_Regression': GD_LR,
          'KNN': GD_KNN,
          'Decision_Trees': GD_DT,
          'Extra_Trees': GD_ET,
          'Random_Forest': GD_RF,
          'Gradient_Boosting': GD_GB}

Avoid the warnings

In [44]:
import warnings
warnings.filterwarnings('ignore')

Train the models

In [45]:
results_col = ['Classifier', 'Accuracy', 'Precision', 'Recall', 'ROC-AUC', 'F1-score']
results = pd.DataFrame(columns=results_col)

l = 0

accuracy_models_opt, precision_models_opt, recall_models_opt = [], [], []
roc_auc_models_opt, f1_score_models_opt = [], []
best_estimator = {}

for name, model in models.items():
    print(name)
    model.fit(X_train,y_train)

    accuracy = cross_val_score(model.best_estimator_, X, y, cv=10, scoring='accuracy')
    accuracy_models_opt.append(accuracy)

    precision = cross_val_score(model.best_estimator_, X, y, cv=10, scoring='precision')
    precision_models_opt.append(precision)

    recall = cross_val_score(model.best_estimator_, X, y, cv=10, scoring='recall')
    recall_models_opt.append(recall)

    roc_auc = cross_val_score(model.best_estimator_, X, y, cv=10, scoring='roc_auc')
    roc_auc_models_opt.append(roc_auc)

    f1 = cross_val_score(model.best_estimator_, X, y, cv=10, scoring='f1')
    f1_score_models_opt.append(f1_score_models_opt)

    best_estimator[name] = model.best_estimator_

    results.loc[l] = [name,
                     '%s %pm2 %s' % (round(accuracy.mean()*100,2),round(accuracy.std()*100,2)),
                     '%s %pm2 %s' % (round(precision.mean()*100,2),round(precision.std()*100,2)),
                     '%s %pm2 %s' % (round(recall.mean()*100,2),round(recall.std()*100,2)),
                     '%s %pm2 %s' % (round(roc_auc.mean()*100,2),round(roc_auc.std()*100,2)),
                     '%s %pm2 %s' % (round(f1.mean()*100,2),round(f1.std()*100,2))]

    l += 1

Logistic_Regression
KNN
Decision_Trees
Extra_Trees
Random_Forest
Gradient_Boosting

Let's see the results

In [46]:
results

Out[46]:
   Classifier  Accuracy  Precision  Recall  ROC-AUC  F1-score
0  Logistic_Regression  91.43 ± 0.58  47.54 ± 14.86  18.95 ± 7.96  82.67 ± 6.55  26.46 ± 9.63
1      KNN              91.6 ± 1.25  50.43 ± 6.86  44.05 ± 11.34  85.49 ± 5.53  46.1 ± 7.99
2  Decision_Trees      89.45 ± 1.99  41.85 ± 7.83  65.05 ± 10.21  79.16 ± 4.64  50.78 ± 6.45
3  Extra_Trees          90.09 ± 1.86  43.17 ± 4.12  64.95 ± 10.94  83.45 ± 6.2  53.72 ± 8.73
4  Random_Forest        89.74 ± 1.97  41.97 ± 6.79  65.05 ± 9.16  88.1 ± 3.47  52.63 ± 7.67
5  Gradient_Boosting     92.54 ± 1.16  62.13 ± 19.05  30.71 ± 9.22  87.66 ± 4.94  40.35 ± 9.91

Let's plot the boxplots of eac classifier and for each scoring metrics

In [47]:
fig, ax = plt.subplots(5, 1, figsize=(20,11))
plt.subplots_adjust(wspace=0.4, hspace=0.4)

ax[0].boxplot(accuracy_models_opt)
ax[1].boxplot(precision_models_opt)
ax[2].boxplot(recall_models_opt)
ax[3].boxplot(roc_auc_models_opt)
ax[4].boxplot(f1_score_models_opt)

ax[0].set_ylabel('Accuracy', labelpad=20, fontsize=16)
ax[1].set_ylabel('Precision', labelpad=20, fontsize=16)
ax[2].set_ylabel('Recall', labelpad=20, fontsize=16)
ax[3].set_ylabel('ROC-AUC', labelpad=20, fontsize=16)
ax[4].set_ylabel('F1-score', labelpad=20, fontsize=16)

xticks = np.arange(1,7)
for i in range(len(ax)):
    if i==4:
        ax[i].set_xticks(xticks)
        ax[i].set_xticklabels(models.keys(), fontsize=17, rotation=75)
    else:
        ax[i].set_xticks(xticks)
        ax[i].set_xticklabels('')

Accuracy
Precision
Recall
ROC-AUC
F1-score

I select 4 best models based on the highest result of the respective metrics. If a model is already included, I choose the next model.

In [48]:
metrics_used = ['Accuracy', 'Precision', 'Recall', 'F1-score']
specific_models = []

for x in metrics_used:
    which_algorithm = results.sort_values(x, ascending=False).iloc[0,0]
    if x != 'Accuracy':
        while which_algorithm in specific_models:
            which_algorithm = results.sort_values(x, ascending=False).iloc[1,0]
            if which_algorithm in specific_models:
                i+=1
            else:
                specific_models.append(which_algorithm)
                break
        else:
            specific_models.append(which_algorithm)
    else:
        specific_models.append(which_algorithm)

specific_models

Out[48]:
['Gradient_Boosting', 'KNN', 'Random_Forest', 'Extra_Trees']

In [49]:
fig, ax = plt.subplots(2, 4, figsize=(18,8))
plt.subplots_adjust(wspace=0.4, hspace=0.4)

# Confusion Matrix
# ROC curve

l = 0
for name in specific_models:
    y_pred = models[name].best_estimator_.predict(X_test)
    cm = confusion_matrix(y_test,y_pred)

    sns.heatmap(cm, annot=True, annot_kws={"size": 10, "weight": "bold", "fmt": 'g', ax=ax[0,l]})

    ticks = ['Pass', 'Fail']
    ax[0,l].set_yticklabels(ticks, fontsize=13)
    ax[0,l].set_xticklabels(ticks, fontsize=13)
    ax[0,l].set_xlabel('Actual value', labelpad=12, fontsize=14)
    ax[0,l].set_ylabel('Predicted value', labelpad=12, fontsize=14)
    ax[0,l].set_title(specific_models[i], y=1.02, fontsize=15)

    i+=1

# ROC curve

l = 0
for i in range(len(specific_models)):
    name = specific_models[i]
    probs = models[name].best_estimator_.predict_proba(X_test)
    probs = probs[:, 1]
    fpr, tpr, _ = roc_curve(y_test, probs)

    ax[1,l].plot(fpr, tpr, lw=2, color='r')
    ax[1,l].plot([0,0.2], [0,0.2], lw=1, color='b')

    ticks = [0, 0.2, 0.4, 0.6, 0.8, 1.0]
    ax[1,l].set_xticks(ticks)
    ax[1,l].set_yticks(ticks, fontsize=13)

    ax[1,l].set_xlim([0,0.2], [0.2,1.025])
    ax[1,l].set_ylim([0,0.2], [0.2,1.025])

    ax[1,l].set_xlabel('False Positive Rate', fontsize=14, labelpad=10)
    ax[1,l].set_ylabel('True Positive Rate', labelpad=10)

    i+=1

Gradient_Boosting
KNN
Random_Forest
Extra_Trees

Feature Importances

In [50]:
fig, ax = plt.subplots(2,2,figsize=(20,17))
plt.subplots_adjust(wspace=0.4, hspace=0.4)
ax = ax.flatten()

for i, name in enumerate(['Extra_Trees', 'Random_Forest', 'Gradient_Boosting']):
    best_model = models[name].best_estimator_
    feature_importances = pd.Series(best_model.feature_importances_,
                                   data_for_predictions.columns[1:-1]).sort_values(ascending=True)

    feature_importances.plot(kind='bar', fontsize=15, ax=axax[i])
    ax[i].set_xlabel('Feature', labelpad=15, fontsize=20)
    ax[i].set_ylabel('Relative Importances', labelpad=15, fontsize=20)
    ax[i].set_title(name, y=1.02, fontsize=20)

Extra_Trees
Random_Forest

Relative Importances
Feature

Gradient_Boosting
Relative Importances
Feature

In [51]:
fig, ax = plt.subplots(1,2,figsize=(20,6))

for i, name in enumerate(['Extra_Trees', 'Random_Forest']):
    best_model = models[name].best_estimator_
    feature_importances = pd.Series(best_model.feature_importances_,
                                   data_for_predictions.columns[1:-1]).sort_values(ascending=True)

    feature_importances.plot(kind='bar', fontsize=15, ax=axax[i])
    ax[i].set_xlabel('Feature', labelpad=15, fontsize=20)
    ax[i].set_ylabel('Relative Importances', labelpad=15, fontsize=20)
    ax[i].set_title(name, y=1.02, fontsize=20)

Extra_Trees
Random_Forest

Relative Importances
Feature

Let's find out the optimal model based on the F1-score:

In [52]:
which_algorithm = results[results['F1-score']] = results['F1-score'].max() ['Classifier']
best_model_method_1 = best_estimator[which_algorithm.values[0]]

results[results['F1-score']] = results['F1-score'].max()

Out[52]:
   Classifier  Accuracy  Precision  Recall  ROC-AUC  F1-score
3  Extra_Trees  90.09 ± 1.86  43.17 ± 4.12  64.95 ± 10.94  83.45 ± 6.2  53.72 ± 8.73

The best performed classifier had taken a parameter which is helpful to balance the weights when the target variable is imbalanced, the class_weight='balanced'. Before I store the predictions in a new column of the dataset, I would like to check how would be result if this parameter is missing...

In [53]:
parameters = {'max_leaf_nodes': [i for i in range(2,10)],
              'max_depth': [i for i in range(2,10)]}

#Decision Trees
GD_DT = GridSearchCV(DecisionTreeClassifier(),
                    param_grid=parameters,
                    scoring='f1')

#Extra Trees
GD_ET = GridSearchCV(ExtraTreeClassifier(),
                    param_grid=parameters,
                    scoring='f1')

#Random Forest
GD_RF = GridSearchCV(RandomForestClassifier(),
                    param_grid=parameters,
                    scoring='f1')

#Gradient Boosting
GD_GB = GridSearchCV(GradientBoostingClassifier(),
                    param_grid=parameters,
                    scoring='f1')

#K-NN
parameters = {'n_neighbors': [i for i in range(5,35,5)]}
GD_KNN = GridSearchCV(KNeighborsClassifier(),
                    param_grid=parameters,
                    scoring='f1')

models = {'Logistic_Regression': GD_LR,
          'KNN': GD_KNN,
          'Decision_Trees': GD_DT,
          'Extra_Trees': GD_ET,
          'Random_Forest': GD_RF,
          'Gradient_Boosting': GD_GB}

Store the predictions into a new column. Firstly, copy the data into a new variable in order to store the predictions.

In [54]:
data_method_predictions = data_for_predictions.copy()

In [55]:
predictions_without_weights = model_without_weights.predict(X)
data_method_predictions['Pred. without weights'] = predictions_without_weights

predictions_method_1 = best_model_method_1.predict(X)
data_method_predictions['Pred. Method 1'] = predictions_method_1

Let's see the classification report for the best model to whole dataset.

In [56]:
from sklearn.metrics import classification_report

print('Using the principle of method 1')
print()
print(classification_report(data_method_predictions.Class,data_method_predictions['Pred. Method 1']))
print('='*60)
print()
print('Best estimator of method 1 without weights')
print()
print(classification_report(data_method_predictions.Class,data_method_predictions['Pred. without weights']))

Using the principle of method 1

           precision    recall  f1-score   support

      0       0.97       0.92       0.95       1572
      1       0.43       0.49       0.46       143

 accuracy: 0.71
 macro avg: 0.71
 weighted avg: 0.93

Best estimator of method 1 without weights

           precision    recall  f1-score   support

      0       0.92       1.00       0.96       1572
      1       0.75       0.04       0.08       143

 accuracy: 0.81
 macro avg: 0.91
 weighted avg: 0.93

# ## Method 2: Employing an experiment by resizing the classes of the target variable

In [57]:
parameters = {'penalty':['l1','l2'],
              'C': [0.001, 0.01, 0.1, 1]}

# Logistic Regression
GD_LR = GridSearchCV(LogisticRegression(solver='liblinear'),
                    param_grid=parameters,
                    scoring='f1')

#Decision Trees
parameters = {'max_leaf_nodes': [i for i in range(2,10)],
              'max_depth': [i for i in range(2,10)]}

GD_DT = GridSearchCV(DecisionTreeClassifier(),
                    param_grid=parameters,
                    scoring='f1')

#Extra Trees
GD_ET = GridSearchCV(ExtraTreeClassifier(),
                    param_grid=parameters,
                    scoring='f1')

#Random Forest
GD_RF = GridSearchCV(RandomForestClassifier(),
                    param_grid=parameters,
                    scoring='f1')

#Gradient Boosting
GD_GB = GridSearchCV(GradientBoostingClassifier(),
                    param_grid=parameters,
                    scoring='f1')

#K-NN
parameters = {'n_neighbors': [i for i in range(5,35,5)]}
GD_KNN = GridSearchCV(KNeighborsClassifier(),
                    param_grid=parameters,
                    scoring='f1')

models = {'Logistic_Regression': GD_LR,
          'KNN': GD_KNN,
          'Decision_Trees': GD_DT,
          'Extra_Trees': GD_ET,
          'Random_Forest': GD_RF,
          'Gradient_Boosting': GD_GB}

import random

size = [1,1.5,2,2.5,3,3.5,5,7.5,10]

accuracy_per_size, precision_per_size, recall_per_size = [], [], []
roc_auc_per_size, f1_per_size = [], [], []
best_estimator_per_size = {}

experiment, classes_size = [], []

print('='*35)
for s in size:
    print('Class Pass -> ', s, ' 1 <- Class Fail')
    print()
    ind_class_1 = data_for_predictions[data_for_predictions.Class==1].index.tolist()
    ind_class_0 = random.choices(data_for_predictions.Class==0).index, k=int(s*len(ind_class_1))
    which_ind = ind_class_1 + ind_class_0
    random.shuffle(which_ind)

    X, y = data_for_predictions.iloc[which_ind,:], data_for_predictions.iloc[which_ind,-1]
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=42)

    classes_size.append(round((y.value_counts()['y'].count()*100,2),values))

    results_col = ['Classifier', 'Accuracy', 'Precision', 'Recall', 'ROC-AUC', 'F1-score']
    results = pd.DataFrame(columns=results_col)

    l = 0

    accuracy_models_opt, precision_models_opt, recall_models_opt = [], [], []
    roc_auc_models_opt, f1_score_models_opt = [], []
    best_estimator = {}

    for name, model in models.items():
        print(name)
        model.fit(X_train,y_train)

        accuracy = cross_val_score(model.best_estimator_, X, y, cv=10, scoring='accuracy')
        accuracy_models_opt.append(accuracy)

        precision = cross_val_score(model.best_estimator_, X, y, cv=10, scoring='precision')
        precision_models_opt.append(precision)

        recall = cross_val_score(model.best_estimator_, X, y, cv=10, scoring='recall')
        recall_models_opt.append(recall)

        roc_auc = cross_val_score(model.best_estimator_, X, y, cv=10, scoring='roc_auc')
        roc_auc_models_opt.append(roc_auc)

        f1 = cross_val_score(model.best_estimator_, X, y, cv=10, scoring='f1')
        f1_score_models_opt.append(f1_score_models_opt)

        best_estimator[name] = model.best_estimator_

        results.loc[l] = [name,
                         '%s %pm2 %s' % (round(accuracy.mean()*100,2),round(accuracy.std()*100,2)),
                         '%s %pm2 %s' % (round(precision.mean()*100,2),round(precision.std()*100,2)),
                         '%s %pm2 %s' % (round(recall.mean()*100,2),round(recall.std()*100,2)),
                         '%s %pm2 %s' % (round(roc_auc.mean()*100,2),round(roc_auc.std()*100,2)),
                         '%s %pm2 %s' % (round(f1.mean()*100,2),round(f1.std()*100,2))]

        l += 1

    experiment.append(results)

    best_estimator_per_size.append(best_estimator)

    accuracy_per_size.append(accuracy_models_opt)
    precision_per_size.append(precision_models_opt)
    recall_per_size.append(recall_models_opt)
    roc_auc_per_size.append(roc_auc_models_opt)
    f1_per_size.append(f1_score_models_opt)

    print('='*35)

Class Pass -> 1 : 1 <- Class Fail

Logistic_Regression
KNN
Decision_Trees
Extra_Trees
Random_Forest
Gradient_Boosting

Class Pass -> 1.5 : 1 <- Class Fail

Logistic_Regression
KNN
Decision_Trees
Extra_Trees
Random_Forest
Gradient_Boosting

Class Pass -> 2 : 1 <- Class Fail

Logistic_Regression
KNN
Decision_Trees
Extra_Trees
Random_Forest
Gradient_Boosting

Class Pass -> 2.5 : 1 <- Class Fail

Logistic_Regression
KNN
Decision_Trees
Extra_Trees
Random_Forest
Gradient_Boosting

Class Pass -> 3 : 1 <- Class Fail

Logistic_Regression
KNN
Decision_Trees
Extra_Trees
Random_Forest
Gradient_Boosting

Class Pass -> 3.5 : 1 <- Class Fail

Logistic_Regression
KNN
Decision_Trees
Extra_Trees
Random_Forest
Gradient_Boosting

Class Pass -> 5 : 1 <- Class Fail

Logistic_Regression
KNN
Decision_Trees
Extra_Trees
Random_Forest
Gradient_Boosting

Class Pass -> 7.5 : 1 <- Class Fail

Logistic_Regression
KNN
Decision_Trees
Extra_Trees
Random_Forest
Gradient_Boosting

Class Pass -> 10 : 1 <- Class Fail

Logistic_Regression
KNN
Decision_Trees
Extra_Trees
Random_Forest
Gradient_Boosting

In [59]:
l = 0
for name, in models.items():
    results_per_classifier = []
    for i in range(len(experiment)):
        results_per_classifier.append([size[i],classes_size[i][0],classes_size[i][1],
                                     round(accuracy_per_size[i][j].mean()*100,2),
                                     round(precision_per_size[i][j].mean()*100,2),
                                     round(recall_per_size[i][j].mean()*100,2),
                                     round(roc_auc_per_size[i][j].mean()*100,2),
                                     round(f1_per_size[i][j].mean()*100,2)])

    globals()[f'{name}_results'] = pd.DataFrame(results_per_classifier, columns=['Ratio Pass:Fail', 'Percentage Fail', 'Accuracy', 'Precision', 'Recall', 'ROC-AUC', 'F1-score'])

    j+=1

In [60]:
algorithms_results = {'Logistic_Regression': Logistic_Regression_results,
                      'KNN': KNN_results,
                      'Decision_Trees': Decision_Trees_results,
                      'Extra_Trees': Extra_Trees_results,
                      'Random_Forest': Random_Forest_results,
                      'Gradient_Boosting': Gradient_Boosting_results}

In [61]:
fig, ax = plt.subplots(3,2,figsize=(16,14))
plt.subplots_adjust(wspace=0.35, hspace=0.55)
ax[2,1].axis('off')

for algorithm_name, results in algorithms_results.items():
    results.plot(x='Ratio Pass:Fail', y='Accuracy', ax=axax[0,0], label=algorithm_name)
    results.plot(x='Ratio Pass:Fail', y='Precision', ax=axax[0,1], label=algorithm_name)
    results.plot(x='Ratio Pass:Fail', y='Recall', ax=axax[1,1], label=algorithm_name)
    results.plot(x='Ratio Pass:Fail', y='F1-score', ax=axax[2,1], label=algorithm_name)

ax[0,0].set_title('Accuracy', y=1.02, fontsize=18)
ax[1,0].set_title('Precision', y=1.02, fontsize=18)
ax[2,0].set_title('Recall', y=1.02, fontsize=18)
ax[0,1].set_title('F1-score', y=1.02, fontsize=18)

x_ticks = [x for x in range(1,11)]
for i in range(3):
    for j in range(2):
        ax[i,j].set_xticks(x_ticks)
        ax[i,j].set_xticklabels(x_ticks, fontsize=14)
        ax[i,j].set_ylabel('Ratio Pass:Fail', labelpad=15, fontsize=14)
        ax[i,j].set_ylabel('Mean value per experiment', labelpad=15, fontsize=14)
        y_ticks = ax[i,j].get_y_ticks()
        ax[i,j].set_yticklabels(y_ticks, fontsize=14)
        ax[i,j].set_yticklabels('')

ax[0,1].legend('')
ax[1,1].legend('')
ax[2,1].legend('')

Accuracy
ROC-AUC
Precision
Recall
F1-score
```



```
fig, ax = plt.subplots(2,2,figsize=(16,10))
plt.subplots_adjust(wspace=.35,hspace=.55)

for algorithm_name, results in algorithms_results.items():
    results.plot(ax='Ratio Pass:Fail', y='Accuracy', axmax[0,0], label=algorithm_name)
    results.plot(ax='Ratio Pass:Fail', y='Precision', axmax[0,1], label=algorithm_name)
    results.plot(ax='Ratio Pass:Fail', y='Recall', axmax[1,0], label=algorithm_name)
    results.plot(ax='Ratio Pass:Fail', y='F1-score', axmax[1,1], label=algorithm_name)

ax[0,0].set_title('Accuracy', y=1.02, fontsize=18)
ax[0,1].set_title('Precision', y=1.02, fontsize=18)
ax[1,0].set_title('Recall', y=1.02, fontsize=18)
ax[1,1].set_title('F1-score', y=1.02, fontsize=18)

x_ticks = [x for x in range(1,11)]
for i in range(2):
    for j in range(2):
        ax[i,j].set_xticks(x_ticks)
        ax[i,j].set_xticklabels(x_ticks, fontsize=14)
        ax[i,j].set_xlabel('Ratio Pass:Fail', labelpad=15, fontsize=14)
        ax[i,j].set_ylabel('Mean value per experiment', labelpad=15, fontsize=14)

        y_ticks = ax[i,j].get_yticks()
        ax[i,j].set_yticks(y_ticks)
        ax[i,j].set_yticklabels(y_ticks, fontsize=14)

ax[0,1].legend();
ax[1,0].legend();
ax[1,1].legend();
```

The figure displays four line plots arranged in a 2x2 grid, showing the performance of five machine learning models across different 'Ratio Pass:Fail' values (1 to 10). The models are: Logistic Regression (blue), ANN (orange), Decision Trees (green), Extra Trees (purple), and Gradient Boosting (brown). The metrics are: Accuracy, Precision, Recall, and F1-score. The y-axis for all plots is 'Mean value per experiment'.

- Accuracy:** All models show an upward trend as the ratio increases, with Gradient Boosting and Extra Trees generally performing best.
- Precision:** Performance is more varied, with some models showing a peak around a ratio of 3-5 before declining.
- Recall:** Most models show a downward trend as the ratio increases, with Gradient Boosting maintaining higher recall values.
- F1-score:** Similar to precision, F1-scores are more varied across the ratios, with Gradient Boosting and Extra Trees often showing higher values.

```
fig, ax = plt.subplots(9,2,figsize=(20,35))
plt.subplots_adjust(hspace=0.7, wspace=1)

for i in range(len(size)):
    rf = best_estimator_per_size[i][4]
    gb = best_estimator_per_size[i][5]

    best_model = [rf,gb]
    for j in range(2):
        feature_importances = pd.Series(best_model[j].feature_importances_,
                                         data_for_predictions.columns[1:]).sort_values(ascending=True)

        feature_importances.plot(kind='bar', fontsize=10, ax=ax[i,j])

    if j==0:
        ax[i,j].set_ylabel('Pass->rs : 1<-Fail' +size[i]), labelpad=10, fontsize=15)
ax[0,0].set_title('Random Forest', y=1.02, fontsize=15)
ax[0,1].set_title('Gradient Boosting', y=1.02, fontsize=15);
```

The figure displays a 9x2 grid of bar charts showing feature importances for Random Forest and Gradient Boosting models across different size categories (Pass->rs : 1<-Fail). The x-axis for each chart lists the features, and the y-axis shows the importance score. The features are: Name\_1, Name\_2, Name\_3, Name\_4, Name\_5, Name\_6, Name\_7, Name\_8, Name\_9, Name\_10, Name\_11, Name\_12, Name\_13, Name\_14, Name\_15, Name\_16, Name\_17, Name\_18, Name\_19, Name\_20, Name\_21, Name\_22, Name\_23, Name\_24, Name\_25, Name\_26, Name\_27, Name\_28, Name\_29, Name\_30, Name\_31, Name\_32, Name\_33, Name\_34, Name\_35, Name\_36, Name\_37, Name\_38, Name\_39, Name\_40, Name\_41, Name\_42, Name\_43, Name\_44, Name\_45, Name\_46, Name\_47, Name\_48, Name\_49, Name\_50, Name\_51, Name\_52, Name\_53, Name\_54, Name\_55, Name\_56, Name\_57, Name\_58, Name\_59, Name\_60, Name\_61, Name\_62, Name\_63, Name\_64, Name\_65, Name\_66, Name\_67, Name\_68, Name\_69, Name\_70, Name\_71, Name\_72, Name\_73, Name\_74, Name\_75, Name\_76, Name\_77, Name\_78, Name\_79, Name\_80, Name\_81, Name\_82, Name\_83, Name\_84, Name\_85, Name\_86, Name\_87, Name\_88, Name\_89, Name\_90, Name\_91, Name\_92, Name\_93, Name\_94, Name\_95, Name\_96, Name\_97, Name\_98, Name\_99, Name\_100, Name\_101, Name\_102, Name\_103, Name\_104, Name\_105, Name\_106, Name\_107, Name\_108, Name\_109, Name\_110, Name\_111, Name\_112, Name\_113, Name\_114, Name\_115, Name\_116, Name\_117, Name\_118, Name\_119, Name\_120, Name\_121, Name\_122, Name\_123, Name\_124, Name\_125, Name\_126, Name\_127, Name\_128, Name\_129, Name\_130, Name\_131, Name\_132, Name\_133, Name\_134, Name\_135, Name\_136, Name\_137, Name\_138, Name\_139, Name\_140, Name\_141, Name\_142, Name\_143, Name\_144, Name\_145, Name\_146, Name\_147, Name\_148, Name\_149, Name\_150, Name\_151, Name\_152, Name\_153, Name\_154, Name\_155, Name\_156, Name\_157, Name\_158, Name\_159, Name\_160, Name\_161, Name\_162, Name\_163, Name\_164, Name\_165, Name\_166, Name\_167, Name\_168, Name\_169, Name\_170, Name\_171, Name\_172, Name\_173, Name\_174, Name\_175, Name\_176, Name\_177, Name\_178, Name\_179, Name\_180, Name\_181, Name\_182, Name\_183, Name\_184, Name\_185, Name\_186, Name\_187, Name\_188, Name\_189, Name\_190, Name\_191, Name\_192, Name\_193, Name\_194, Name\_195, Name\_196, Name\_197, Name\_198, Name\_199, Name\_200, Name\_201, Name\_202, Name\_203, Name\_204, Name\_205, Name\_206, Name\_207, Name\_208, Name\_209, Name\_210, Name\_211, Name\_212, Name\_213, Name\_214, Name\_215, Name\_216, Name\_217, Name\_218, Name\_219, Name\_220, Name\_221, Name\_222, Name\_223, Name\_224, Name\_225, Name\_226, Name\_227, Name\_228, Name\_229, Name\_230, Name\_231, Name\_232, Name\_233, Name\_234, Name\_235, Name\_236, Name\_237, Name\_238, Name\_239, Name\_240, Name\_241, Name\_242, Name\_243, Name\_244, Name\_245, Name\_246, Name\_247, Name\_248, Name\_249, Name\_250, Name\_251, Name\_252, Name\_253, Name\_254, Name\_255, Name\_256, Name\_257, Name\_258, Name\_259, Name\_260, Name\_261, Name\_262, Name\_263, Name\_264, Name\_265, Name\_266, Name\_267, Name\_268, Name\_269, Name\_270, Name\_271, Name\_272, Name\_273, Name\_274, Name\_275, Name\_276, Name\_277, Name\_278, Name\_279, Name\_280, Name\_281, Name\_282, Name\_283, Name\_284, Name\_285, Name\_286, Name\_287, Name\_288, Name\_289, Name\_290, Name\_291, Name\_292, Name\_293, Name\_294, Name\_295, Name\_296, Name\_297, Name\_298, Name\_299, Name\_300, Name\_301, Name\_302, Name\_303, Name\_304, Name\_305, Name\_306, Name\_307, Name\_308, Name\_309, Name\_310, Name\_311, Name\_312, Name\_313, Name\_314, Name\_315, Name\_316, Name\_317, Name\_318, Name\_319, Name\_320, Name\_321, Name\_322, Name\_323, Name\_324, Name\_325, Name\_326, Name\_327, Name\_328, Name\_329, Name\_330, Name\_331, Name\_332, Name\_333, Name\_334, Name\_335, Name\_336, Name\_337, Name\_338, Name\_339, Name\_340, Name\_341, Name\_342, Name\_343, Name\_344, Name\_345, Name\_346, Name\_347, Name\_348, Name\_349, Name\_350, Name\_351, Name\_352, Name\_353, Name\_354, Name\_355, Name\_356, Name\_357, Name\_358, Name\_359, Name\_360, Name\_361, Name\_362, Name\_363, Name\_364, Name\_365, Name\_366, Name\_367, Name\_368, Name\_369, Name\_370, Name\_371, Name\_372, Name\_373, Name\_374, Name\_375, Name\_376, Name\_377, Name\_378, Name\_379, Name\_380, Name\_381, Name\_382, Name\_383, Name\_384, Name\_385, Name\_386, Name\_387, Name\_388, Name\_389, Name\_390, Name\_391, Name\_392, Name\_393, Name\_394, Name\_395, Name\_396, Name\_397, Name\_398, Name\_399, Name\_400, Name\_401, Name\_402, Name\_403, Name\_404, Name\_405, Name\_406, Name\_407, Name\_408, Name\_409, Name\_410, Name\_411, Name\_412, Name\_413, Name\_414, Name\_415, Name\_416, Name\_417, Name\_418, Name\_419, Name\_420, Name\_421, Name\_422, Name\_423, Name\_424, Name\_425, Name\_426, Name\_427, Name\_428, Name\_429, Name\_430, Name\_431, Name\_432, Name\_433, Name\_434, Name\_435, Name\_436, Name\_437, Name\_438, Name\_439, Name\_440, Name\_441, Name\_442, Name\_443, Name\_444, Name\_445, Name\_446, Name\_447, Name\_448, Name\_449, Name\_450, Name\_451, Name\_452, Name\_453, Name\_454, Name\_455, Name\_456, Name\_457, Name\_458, Name\_459, Name\_460, Name\_461, Name\_462, Name\_463, Name\_464, Name\_465, Name\_466, Name\_467, Name\_468, Name\_469, Name\_470, Name\_471, Name\_472, Name\_473, Name\_474, Name\_475, Name\_476, Name\_477, Name\_478, Name\_479, Name\_480, Name\_481, Name\_482, Name\_483, Name\_484, Name\_485, Name\_486, Name\_487, Name\_488, Name\_489, Name\_490, Name\_491, Name\_492, Name\_493, Name\_494, Name\_495, Name\_496, Name\_497, Name\_498, Name\_499, Name\_500, Name\_501, Name\_502, Name\_503, Name\_504, Name\_505, Name\_506, Name\_507, Name\_508, Name\_509, Name\_510, Name\_511, Name\_512, Name\_513, Name\_514, Name\_515, Name\_516, Name\_517, Name\_518

[illegible]

```

In [79]:
#Decision Trees
parameters = {'max_leaf_nodes': [1 for i in range(2,10)],
              'max_depth': [1 for i in range(2,10)]}

GD_DT = GridSearchCV(DecisionTreeClassifier(class_weight="balanced"),
                    param_grid=parameters,
                    scoring='f1')

#Extra Trees
GD_ET = GridSearchCV(ExtraTreesClassifier(class_weight="balanced"),
                    param_grid=parameters,
                    scoring='f1')

#Random Forest
GD_RF = GridSearchCV(RandomForestClassifier(class_weight="balanced"),
                    param_grid=parameters,
                    scoring='f1')

#Gradient Boosting
GD_GB = GridSearchCV(GradientBoostingClassifier(),
                    param_grid=parameters,
                    scoring='f1')

#k-NN
parameters = {'n_neighbors': [i for i in range(5,35,5)]}
GD_KNN = GridSearchCV(KNeighborsClassifier(),
                    param_grid=parameters,
                    scoring='f1')

models = ('Logistic_Regression': GD_LR,
          'KNN': GD_KNN,
          'Decision_Trees': GD_DT,
          'Extra_Trees': GD_ET,
          'Random_Forest': GD_RF,
          'Gradient_Boosting': GD_GB)

Train the models

In [79]:
results_col = ['Classifier', 'Accuracy', 'Precision', 'Recall', 'ROC-AUC', 'F1-score']
results = pd.DataFrame(columns=results_col)

l = 0

accuracy_models_opt, precision_models_opt, recall_models_opt = [], [], []
roc_auc_models_opt, f1_score_models_opt = [], []
best_estimator = []

for name, model in models.items():
    print(name)
    model.fit(X_train,y_train)

    accuracy = cross_val_score(model.best_estimator_, X, y, cv=10, scoring="accuracy")
    accuracy_models_opt.append(accuracy)
    precision = cross_val_score(model.best_estimator_, X, y, cv=10, scoring="precision")
    precision_models_opt.append(precision)
    recall = cross_val_score(model.best_estimator_, X, y, cv=10, scoring="recall")
    recall_models_opt.append(recall)
    roc_auc = cross_val_score(model.best_estimator_, X, y, cv=10, scoring="roc_auc")
    roc_auc_models_opt.append(roc_auc)
    f1 = cross_val_score(model.best_estimator_, X, y, cv=10, scoring="f1")
    f1_score_models_opt.append(f1)

    best_estimator[name] = model.best_estimator_

    results.loc[l] = [name,
                     '%s' % ' '.join([round(accuracy.mean()*100,2), round(accuracy.std()*100,2),
                                       '%s' % ' '.join([round(precision.mean()*100,2), round(precision.std()*100,2),
                                       '%s' % ' '.join([round(recall.mean()*100,2), round(recall.std()*100,2),
                                       '%s' % ' '.join([round(roc_auc.mean()*100,2), round(roc_auc.std()*100,2),
                                       '%s' % ' '.join([round(f1.mean()*100,2), round(f1.std()*100,2)])])])])])])

    l += 1

Logistic_Regression
KNN
Decision_Trees
Extra_Trees
Random_Forest
Gradient_Boosting

In [80]:
results

Out[80]:


|   | Classifier          | Accuracy     | Precision     | Recall        | ROC-AUC      | F1-score      |
|---|---------------------|--------------|---------------|---------------|--------------|---------------|
| 0 | Logistic_Regression | 91.78 ± 0.92 | 56.33 ± 23.73 | 14.05 ± 5.61  | 77.72 ± 7.47 | 22.03 ± 8.49  |
| 1 | KNN                 | 92.13 ± 1.48 | 59.08 ± 16.77 | 34.86 ± 13.11 | 85.28 ± 5.14 | 41.47 ± 10.49 |
| 2 | Decision_Trees      | 89.97 ± 1.55 | 42.98 ± 6.61  | 61.48 ± 9.04  | 79.42 ± 4.58 | 50.41 ± 6.97  |
| 3 | Extra_Trees         | 90.14 ± 1.45 | 42.11 ± 6.57  | 65.05 ± 9.7   | 78.77 ± 7.0  | 52.67 ± 7.14  |
| 4 | Random_Forest       | 90.49 ± 1.53 | 44.74 ± 6.0   | 62.9 ± 8.8    | 83.88 ± 5.34 | 52.67 ± 6.7   |
| 5 | Gradient_Boosting   | 92.01 ± 1.57 | 53.98 ± 19.78 | 28.0 ± 13.14  | 87.06 ± 3.02 | 35.73 ± 14.51 |



I select 4 best models based on the highest result of the respective metrics. If a model is already included, I choose the next model.

In [81]:
metrics_used = ['Accuracy', 'Precision', 'Recall', 'F1-score']

specific_models = []

for x in metrics_used:
    which_algorithm = results.sort_values(x, ascending=False).iloc[0,0]
    if x!="Accuracy":
        if which_algorithm in specific_models:
            i=1
            while which_algorithm in specific_models:
                which_algorithm = results.sort_values(x, ascending=False).iloc[i,0]
                if which_algorithm in specific_models:
                    i+=1
                else:
                    specific_models.append(which_algorithm)
            break
        else:
            specific_models.append(which_algorithm)
    else:
        specific_models.append(which_algorithm)

specific_models

Out[81]:
['KNN', 'Logistic_Regression', 'Extra_Trees', 'Random_Forest']

```

```

fig, ax = plt.subplots(2, 4, figsize=(15,20))
plt.subplots_adjust(wspace=0.4, hspace=0.4)

#####
#
# Confusion Matrix
#
#####

i = 0
for name in specific_models:
    y_pred = models[name].best_estimator_.predict(X_test)
    cm = confusion_matrix(y_test, y_pred)

    sns.heatmap(cm, annot=True, annot_kws={"size": 10, "weight": "bold", "fmt": 'g', "axmax[0,1]}

    ticks = ['Pass', 'Fail']
    ax[0,i].set_yticklabels(ticks, fontsize=13)
    ax[0,i].set_xticklabels(ticks, fontsize=13)
    ax[0,i].set_xlabel('Actual value', labelpad=12, fontsize=14)
    ax[0,i].set_ylabel('Predicted value', labelpad=12, fontsize=14)

    ax[0,i].set_title(specific_models[i], y=1.02, fontsize=15)
    i+=1

#####
#
# ROC curve
#
#####

i = 0
for name in specific_models:
    probs = models[name].best_estimator_.predict_proba(X_test)
    probs = probs[:, 1]
    rf_fpr, rf_tpr, rf_thresholds = roc_curve(y_test, probs)

    ax[1,i].plot(rf_fpr, rf_tpr, '-o')
    ax[1,i].plot([-0.02,1],[-0.02,1], '-r')

    ticks = [0, 0.2, 0.4, 0.6, 0.8, 1.0]
    ax[1,i].set_xticks(ticks)
    ax[1,i].set_yticklabels(ticks, fontsize=13)

    ax[1,i].set_yticks(ticks)

    ax[1,i].set_xlim(-0.02,1.025)
    ax[1,i].set_ylim(-0.02,1.025)

    ax[1,i].set_xlabel('False Positive Rate', fontsize=14, labelpad=10)
    ax[1,i].set_ylabel('True Positive Rate', fontsize=14, labelpad=10)

    i+=1

#####
#
# Feature importances
#
#####

fig, ax = plt.subplots(1,3,figsize=(20,6))

for i, name in enumerate(['Extra_Trees', 'Random_Forest', 'Decision_Trees']):
    best_model = models[name].best_estimator_
    feature_importances = pd.Series(best_model.feature_importances_,
                                   data_fa.columns[-1:]).sort_values(ascending=True)

    feature_importances.plot(kind='bar', fontsize=15, ax=ax[i])
    ax[i].set_xlabel('Feature', labelpad=15, fontsize=20)
    ax[i].set_title(name, y=1.02, fontsize=20)
    if i==0:
        ax[i].set_ylabel('Relative Importances', fontsize=20);

#####
#
# Let's find out the optimal model based on the F1-score:
#
#####

which_algorithm = results[results['F1-score']==results['F1-score'].max()]['Classifier']
best_model_method_3 = best_estimator[which_algorithm.values[0]]

results[results['F1-score']==results['F1-score'].max()]

#####
#
# Store the predictions
#
#####

predictions_method_3 = best_model_method_3.predict(X)
data_method_predictions["Pred. Method 3"] = predictions_without_weights

#####
#
# Conclusions
#
#####

Factor Analysis is a dimensionality reduction method which is mainly used when features are categorical. I transformed the first 4 features because those were enough to overtake the threshold of 80%. The best model has been performed by the Decision Trees. Looking at the confusion matrix and f1-score the most predicted values has been misclassified.

#####
#
# Analyzing the methods
#
#####

fig, ax = plt.subplots(1,4, figsize=(20,5))

data_method_predictions.Class.hist(ax=ax[0], alpha=0.3, label='Actual class')
data_method_predictions["Pred. without weights"].hist(ax=ax[0], alpha=0.3, label='Predicted class')
ax[0].set_title('Predictions without weights', y=1.02, fontsize=18)
ax[0].legend()

data_method_predictions.Class.hist(ax=ax[1], alpha=0.3, label='Actual class')
data_method_predictions["Pred. Method 1"].hist(ax=ax[1], alpha=0.3, label='Predicted class')
ax[1].set_title('Predictions Method 1', y=1.02, fontsize=18)
ax[1].legend()

data_method_predictions.Class.hist(ax=ax[2], alpha=0.3, label='Actual class')
data_method_predictions["Pred. Method 2"].hist(ax=ax[2], alpha=0.3, label='Predicted class')
ax[2].set_title('Predictions Method 2', y=1.02, fontsize=18)
ax[2].legend()

data_method_predictions.Class.hist(ax=ax[3], alpha=0.3, label='Actual class')
data_method_predictions["Pred. Method 3"].hist(ax=ax[3], alpha=0.3, label='Predicted class')
ax[3].set_title('Predictions Method 3', y=1.02, fontsize=18)
ax[3].legend()

#####
#
# Predictions without weights
#
#####

#####
#
# Predictions Method 1
#
#####

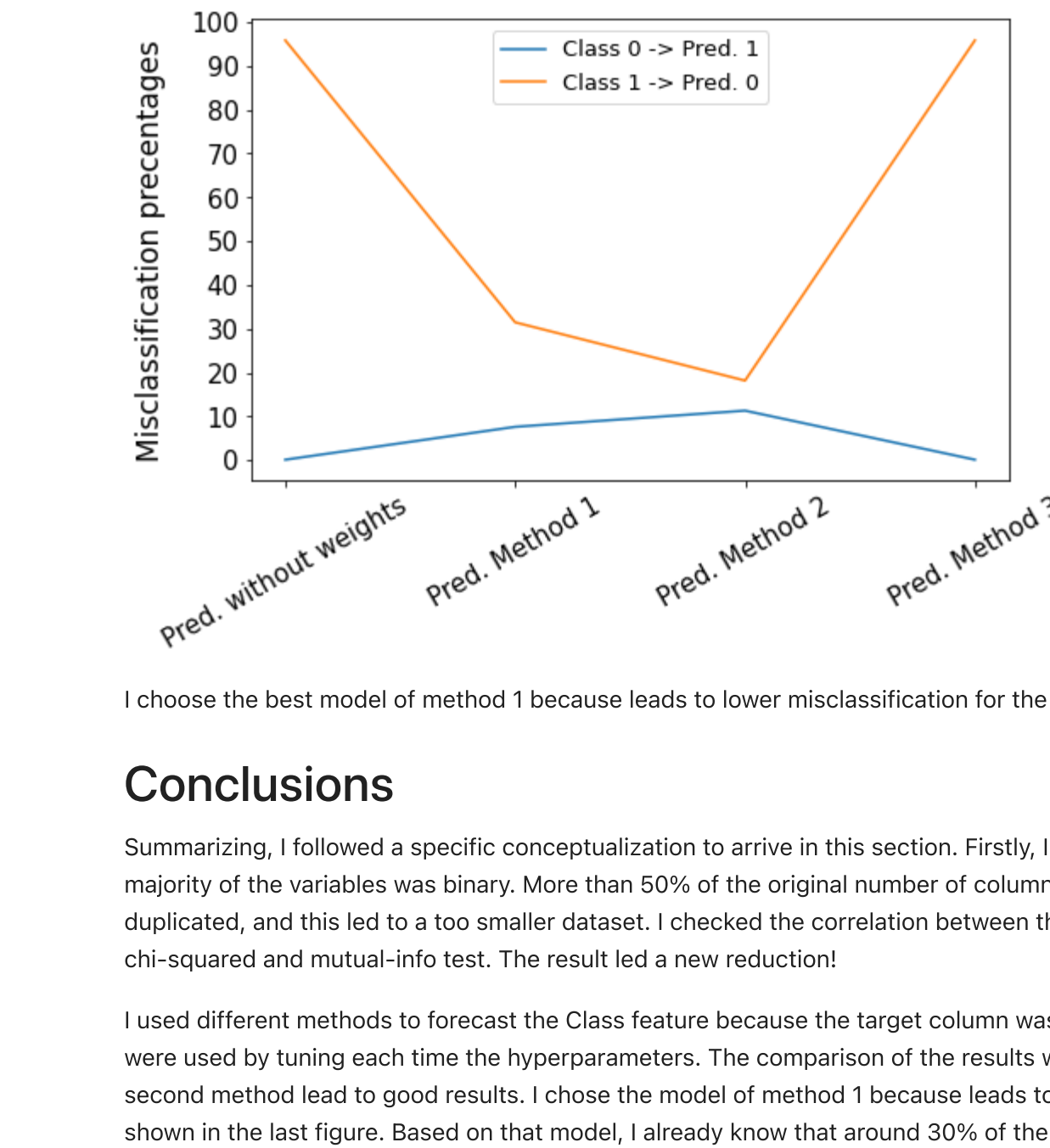
#####
#
# Predictions Method 2
#
#####

#####
#
# Predictions Method 3
#
#####

#####
#
# methods_analysis = data_method_predictions[['Class','Pred. without weights','Pred. Method 1',
#                                           'Pred. Method 2','Pred. Method 3']]
#
#####
#
# columns = ['Pred. Without weights','Pred. Method 1','Pred. Method 2', 'Pred. Method 3']
#
#####
#
# # count how many are in total with 0 and 1 in the variable Class
# how_many_0 = methods_analysis[methods_analysis.Class==0]['Class'].count()
# how_many_1 = methods_analysis[methods_analysis.Class==1]['Class'].count()
#
# perc_mis_0_to_1_per_method, perc_mis_1_to_0_per_method = [], []
#
# for x in columns:
#     mis_0_to_1 = methods_analysis[(methods_analysis.Class==0) &
#                                   (methods_analysis[x]==1)][x].count()
#     mis_1_to_0 = methods_analysis[(methods_analysis.Class==1) &
#                                   (methods_analysis[x]==0)][x].count()
#
#     perc_mis_0_to_1_per_method.append(round(mis_0_to_1/how_many_0*100,2))
#     perc_mis_1_to_0_per_method.append(round(mis_1_to_0/how_many_1*100,2))
#
# correct_classified_0 = methods_analysis[(methods_analysis.Class==0) &
#                                           (methods_analysis[x]==0)][x].count()
# correct_classified_1 = methods_analysis[(methods_analysis.Class==1) &
#                                           (methods_analysis[x]==1)][x].count()
#
# print('For the column', x, ':')
# print()
# print('Classified 0 but predicted to 1:', mis_0_to_1)
# print('Classified 1 but predicted to 0:', mis_1_to_0)
# print('---*40')
#
# For the column Pred. without weights :
#
# Classified 0 but predicted to 1: 2
# Classified 1 but predicted to 0: 137
#
# For the column Pred. Method 1 :
#
# Classified 0 but predicted to 1: 120
# Classified 1 but predicted to 0: 45
#
# For the column Pred. Method 2 :
#
# Classified 0 but predicted to 1: 178
# Classified 1 but predicted to 0: 26
#
# For the column Pred. Method 3 :
#
# Classified 0 but predicted to 1: 2
# Classified 1 but predicted to 0: 137
#
#####
#
# fig, ax = plt.subplots(1,1,figsize=(8,5))
#
#####
#
# pos = np.arange(4)
# ax.plot(pos, perc_mis_0_to_1_per_method, label='Class 0 -> Pred. 1')
# ax.plot(pos, perc_mis_1_to_0_per_method, label='Class 1 -> Pred. 0')
#
# ax.set_xticks(pos)
# ax.set_yticklabels(columns, fontsize=5, rotation=30)
# ax.set_yticks([i for i in range(0,110,10)])
# ax.set_xlabel('Misclassification percentages', labelpad=14, fontsize=15)
#
# ax.set_ylabel('Misclassification percentages', labelpad=14, fontsize=17)
#
# ax.legend(loc='best', fontsize=13);

```





I choose the best model of method 1 because leads to lower misclassification for the minority class.

## Conclusions

Summarizing, I followed a specific conceptualization to arrive in this section. Firstly, I checked for duplicated columns because the majority of the variables was binary. More than 50% of the original number of columns had to be removed, because they were duplicated, and this led to a too smaller dataset. I checked the correlation between the target feature and the binary features using chi-squared and mutual-info test. The result led a new reduction!

I used different methods to forecast the Class feature because the target column was unbalanced. In each method, same algorithms were used by tuning each time the hyperparameters. The comparison of the results was done by using the f1-score. The first and second method lead to good results. I chose the model of method 1 because leads to lower misclassification for the minority class, as shown in the last figure. Based on that model, I already know that around 30% of the predicted class Fail is misclassified and same applies for the predicted class Pass with a percentage of 10%.

Before I move on to the appendix to use the model in order to predict the classes for the Test.csv file, I have to say that the best estimator of the models of methods 1 & 2 needs to run a lot of times for different splitted data by enhancing the hyperparameters and collecting enough statistics to select the best one in the end!

## Appendix - Prepare the submission file

```
In [90]: test_set = pd.read_csv('Test.csv')
test_set.head()
```

	feature_1	feature_2	feature_3	feature_4	feature_5	feature_6	feature_7	feature_8	feature_9	feature_10	...	feature_1549	feature_1550
0	60.0	468.0	78000	1.0	0	0	0	0	0	0	...	0	0
1	108.0	179.0	16574	1.0	0	0	0	0	0	0	...	0	0
2	1.0	1.0	2.0000	0.0	0	0	0	0	0	0	...	0	0
3	60.0	468.0	78000	1.0	0	0	0	0	0	0	...	0	0
4	60.0	120.0	2.0000	1.0	0	0	0	0	0	0	...	0	0

5 rows × 1558 columns

Include the features after the chi\_squared test

```
In [91]: test_set = test_set[continuous_features + data.columns[np.where(scores_chi>=1)].tolist()]
test_set.head()
```

	feature_1	feature_2	feature_3	feature_13	feature_345	feature_954	feature_1044	feature_1048	feature_1110	feature_1137	...	feature_1549	feature_1550
0	60.0	468.0	78000	0	0	0	0	0	0	0	...	0	0
1	108.0	179.0	16574	0	0	0	0	0	0	0	...	1	...
2	1.0	1.0	2.0000	0	0	0	0	0	0	0	...	0	...
3	60.0	468.0	78000	0	0	0	0	0	0	0	...	0	...
4	60.0	120.0	2.0000	0	0	0	0	0	1	0	...	0	...

5 rows × 26 columns

Scale the data

```
In [92]: from sklearn.preprocessing import StandardScaler, MinMaxScaler
ss = StandardScaler()
mm = MinMaxScaler()

for col in test_set.columns[:3]:
    ss.fit(test_set[col].to_numpy().reshape(-1,1))
    test_set[col] = ss.transform(test_set[col].to_numpy().reshape(-1,1))
    test_set[col] = mm.fit_transform(test_set[col].to_numpy().reshape(-1,1)).reshape(-1,1)[0]
```

Apply the model of method 1

I already know that around 28% of Predicted class 1 must have been predicted as 0, and around 8% of class 0 to 1. Those are the misclassifications that I would wait for the test set, if I knew real Class.

```
In [93]: predictions_test_set = best_model_method_1.predict(test_set)
test_set['Predicted Class'] = predictions_test_set
```

```
In [94]: round((test_set['Predicted Class'].value_counts()/test_set.shape[0])*100,2).plot.pie(autopct='%1.2f%%',
figsize=(8,8), fontsize=20)
plt.title('Predicted Class', fontsize=20)
plt.ylabel('');
```

Predicted Class



Submission

```
In [95]: test_set['Predicted Class'].to_csv('Submission.csv', index=False)
```