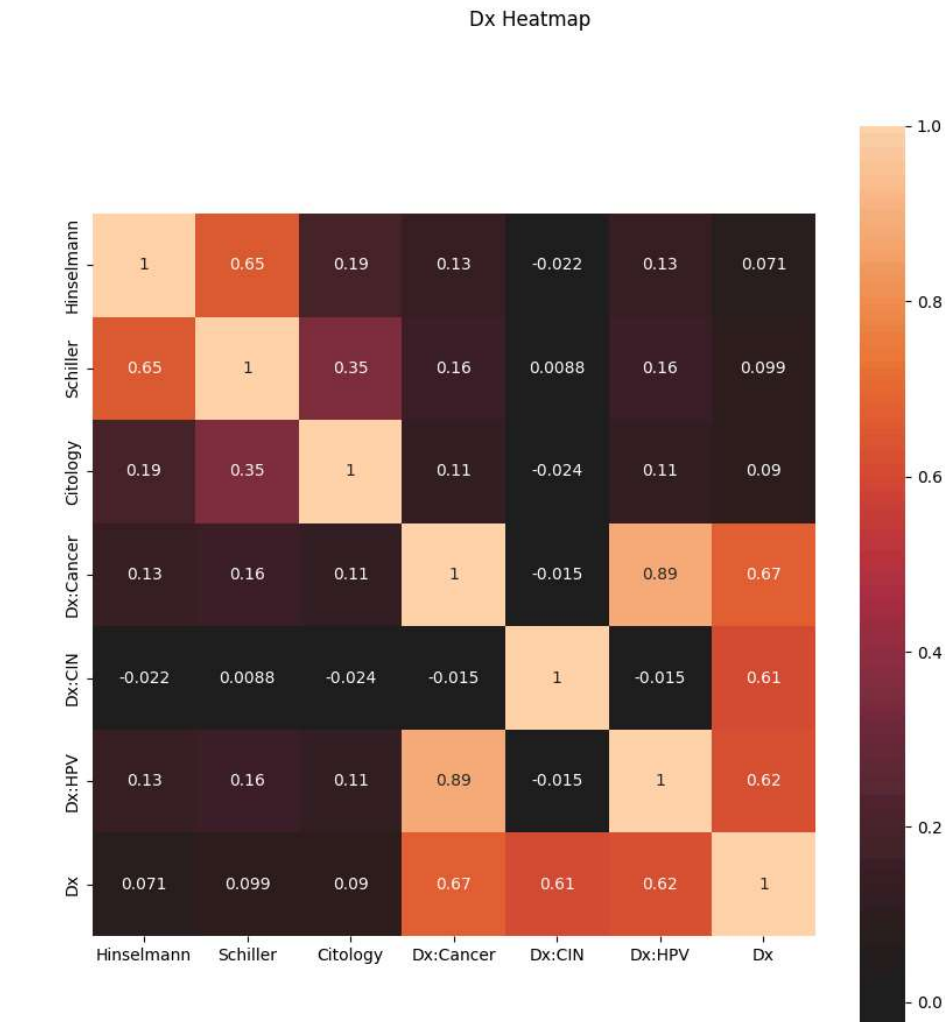


```
plt.figure(figsize=(10,10))
sns.heatmap(
    pd.concat([target_variables, c_dx], axis = 1).corr(),
    annot=True, center=0, square=True
)
plt.suptitle('Dx Heatmap')

Text(0.5, 0.98, 'Dx Heatmap')
```



Logistic Regression

```
# creating the X and Y arrays
X = cervical_cancer.drop(['Hinselmann', 'Schiller', 'Citology', 'Biopsy'], axis=1)
y = cervical_cancer['Hinselmann'] # Only Hinselmann would be used as dependent variable

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)

X_train.shape, X_test.shape

((668, 32), (167, 32))

# Performing Feature Scaling
cols = X_train.columns
scaler = MinMaxScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

X_train = pd.DataFrame(X_train, columns = [cols])
X_test = pd.DataFrame(X_test, columns = [cols])
X_train.describe()
```



```
print('Model accuracy score: {0:0.4f}'.format(accuracy_score(y_test, y_pred_test)))
```

Model accuracy score: 0.9581

[illegible]

```
print('Model accuracy score: {0:0.4f}'.format(accuracy_score(y_train, y_pred_train)))
```

Model accuracy score: 0.9581

```
# Checking for overfitting and underfitting
# getting set scores

print('Training set score: {0:0.4f}'.format(logreg.score(X_train, y_train)))
print('Test set score: {0:0.4f}'.format(logreg.score(X_test, y_test)))
# Since the set score are the same, there is not overfitting

Training set score: 0.9581
Test set score: 0.9581
```

```
# increasing the C to 100
logreg100 = LogisticRegression(C=100, solver='liblinear', random_state=0)
logreg100.fit(X_train, y_train)
```

```
LogisticRegression
LogisticRegression(C=100, random_state=0, solver='liblinear')
```

```
print('Training set score: {:.4f}'.format(logreg100.score(X_train, y_train)))
print('Test set score: {:.4f}'.format(logreg100.score(X_test, y_test)))
# It can be noticed that the set score became higher in the training set but lower on the test set

Training set score: 0.9626
Test set score: 0.9521
```

```
logreg001 = LogisticRegression(C=0.01, solver='liblinear', random_state=0)
logreg001.fit(X_train, y_train)
```

```
LogisticRegression
LogisticRegression(C=0.01, random_state=0, solver='liblinear')
```

```
print('Training set score: {:.4f}'.format(logreg001.score(X_train, y_train)))
print('Test set score: {:.4f}'.format(logreg001.score(X_test, y_test)))
# It can be noticed that there is no difference between a C of 1 and 0.01
```

```

Training set score: 0.9581
Test set score: 0.9581

# Comparing model accuracy with null accuracy
y_test.value_counts()

Hinselmann
0    160
1      7
Name: count, dtype: int64

null_accuracy = (160/167)
print('Null accuracy score: {:.4f}'.format(null_accuracy))
# It can be seen that the null accuracy score is the same as the training set score

Null accuracy score: 0.9581

# Creating an displaying a confusion matrix
cm = confusion_matrix(y_test, y_pred_test)
print('Confusion matrix\n\n', cm)
print('\nTrue Positives(TP) = ', cm[0,0])
print('\nTrue Negatives(TN) = ', cm[1,1])
print('\nFalse Positives(FP) = ', cm[0,1])
print('\nFalse Negatives(FN) = ', cm[1,0])
# This shows that the model created 160 correct predictions and 7 incorrect predictions

Confusion matrix

[[160   0]
 [  7   0]]

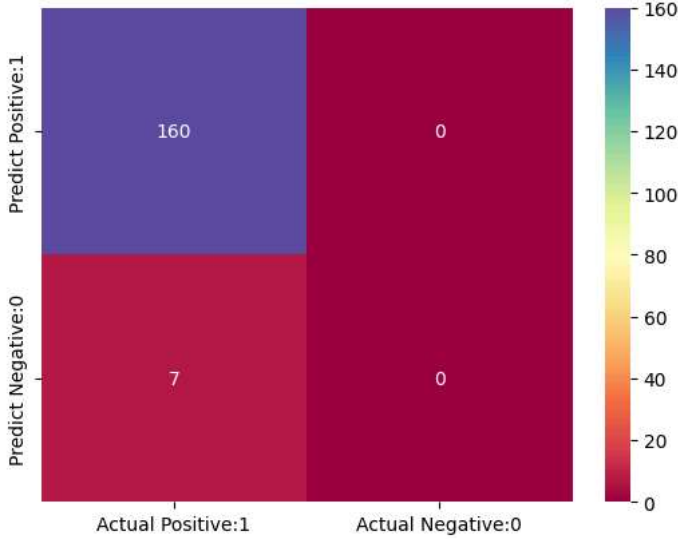
True Positives(TP) = 160

True Negatives(TN) = 0

False Positives(FP) = 0

False Negatives(FN) = 7

# Visualizing the confusion matrix
cm_matrix = pd.DataFrame(data=cm, columns=['Actual Positive:1', 'Actual Negative:0'],
                        index=['Predict Positive:1', 'Predict Negative:0'])
sns.heatmap(cm_matrix, annot=True, fmt='d', cmap='Spectral')

<Axes: >


# Creates a classification report
print(classification_report(y_test, y_pred_test))

              precision    recall  f1-score   support

     0       0.96      1.00      0.98        160
     1       0.00      0.00      0.00         7

 accuracy      0.96      0.96      0.96        167
  macro avg       0.48      0.50      0.49        167
 weighted avg       0.92      0.96      0.94        167

/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning: Precision and F-score are ill-c
_warn_prf(average, modifier, msg_start, len(result))
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning: Precision and F-score are ill-c
_warn_prf(average, modifier, msg_start, len(result))
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning: Precision and F-score are ill-c
_warn_prf(average, modifier, msg_start, len(result))

# Gets the classification accuracy
TP = cm[0,0]
TN = cm[1,1]
FP = cm[0,1]
FN = cm[1,0]

classification_accuracy = (TP + TN) / float(TP + TN + FP + FN)
print('Classification accuracy: {:.4f}'.format(classification_accuracy))

Classification accuracy: 0.9581
```

```
# Gets classification error

classification_error = (FP + FN) / float(TP + TN + FP + FN)
print('Classification error: {:.4f}'.format(classification_error))
# It can be seen that the classification accuracy is high while the error is low

    Classification error: 0.0419

# Gets precision score
precision = TP / float(TP + FP)
print('Precision: {:.4f}'.format(precision))
# The result is a perfect precision score

    Precision: 1.0000

# Gets Recall rate
recall = TP / float(TP + FN)
print('Recall or Sensitivity: {:.4f}'.format(recall))

    Recall or Sensitivity: 0.9581

# Gets False Positive rate
false_positive_rate = FP / float(FP+TN)
print('False Positive: {:.4f}'.format(false_positive_rate))
# Error is encountered as FP is 0

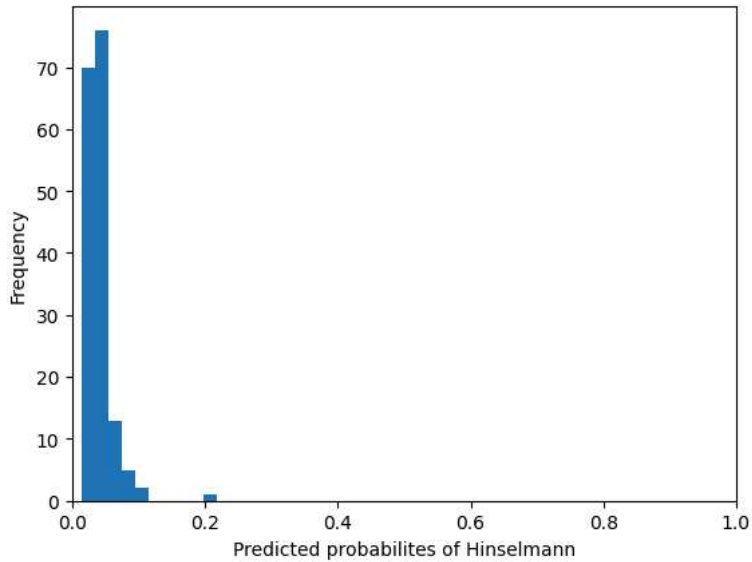
    False Positive: nan
    <ipython-input-118-648a66821cfc>:2: RuntimeWarning: invalid value encountered in divide
        false_positive_rate = FP / float(FP+TN)

# Gets Specificity rate
specificity = TN / (TN + FP)
print('Specificity: {:.4f}'.format(specificity))
# Error is encountered as TN is 0

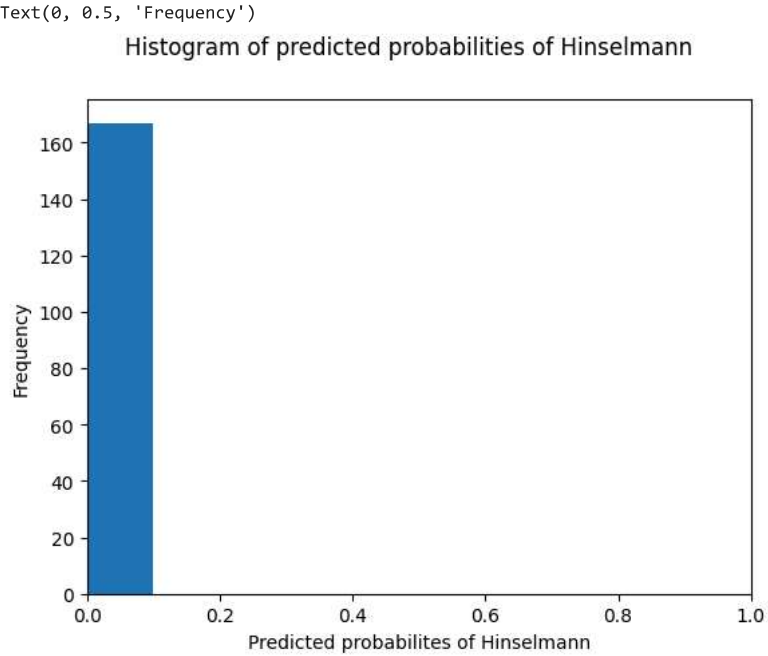
    Specificity: nan
    <ipython-input-119-c484c783968f>:2: RuntimeWarning: invalid value encountered in scalar divide
        specificity = TN / (TN + FP)

# Plots predicted probabilities to a histogram
y_pred1 = logreg.predict_proba(X_test)[: ,1] # using the test variable
plt.hist(y_pred1, bins=10)
plt.suptitle('Histogram of predicted probabilities of Hinselmann')
plt.xlim(0,1)
plt.xlabel('Predicted probabilites of Hinselmann')
plt.ylabel('Frequency')
# The plot is highly skewed to 0,
# which means that the model predicts there a person won't have Hinselmann

    Text(0, 0.5, 'Frequency')
    Histogram of predicted probabilities of Hinselmann
```

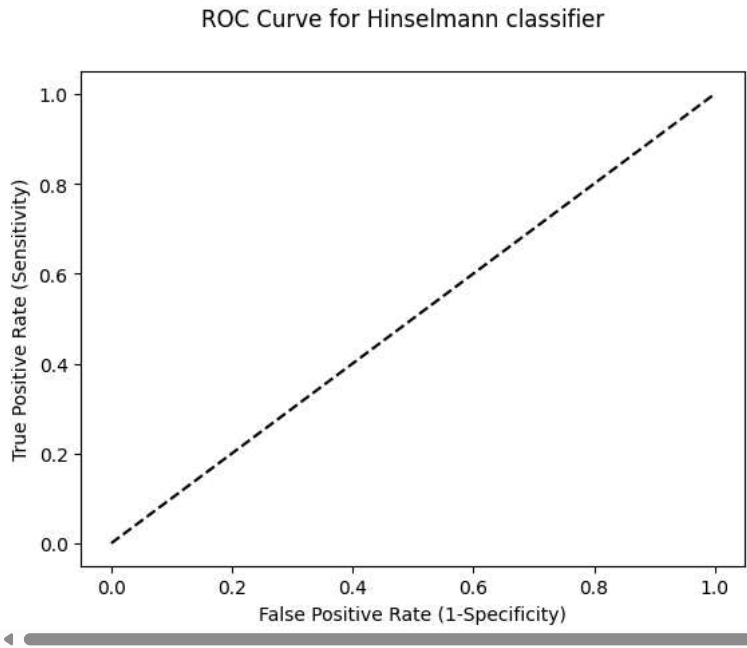


```
# Using 0.5 threshold
y_pred1 = y_pred1.reshape(-1,1)
y_pred2 = binarize(y_pred1, threshold = 0.5)
plt.hist(y_pred2, bins=10)
plt.suptitle('Histogram of predicted probabilities of Hinselmann')
plt.xlim(0,1)
plt.xlabel('Predicted probabilites of Hinselmann')
plt.ylabel('Frequency')
# With the threshold at 0.5, the plot becomes skewed to 0
```



```
# Creates an ROC Curve
fpr, tpr, thresholds = roc_curve(y_train, y_pred_train, pos_label = 'Yes') # Using train variables
plt.plot(fpr, tpr, linewidth = 2)
plt.plot([0,1], [0,1], 'k--')
plt.suptitle('ROC Curve for Hinselmann classifier')
plt.xlabel('False Positive Rate (1-Specificity)')
plt.ylabel('True Positive Rate (Sensitivity)')
# It seems that there is an error as the false positive is 0

/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_ranking.py:1029: UndefinedMetricWarning: No positive samples in y_true, true pc
warnings.warn(
Text(0, 0.5, 'True Positive Rate (Sensitivity)')
```



```
# Gets the ROC AUC score
Cross_validated_ROC_AUC = cross_val_score(logreg, X_train, y_train, cv=5, scoring='roc_auc').mean()
print('Cross validated ROC AUC: {:.4f}'.format(Cross_validated_ROC_AUC))
# The result is not too high but is still higher than 0.5,
# therefore the model might not predict well
```

Cross validated ROC AUC: 0.6114

```
# Gets the 5-Fold Cross Validation score
scores = cross_val_score(logreg, X_train, y_train, cv=5, scoring='accuracy').mean()
print('Cross-validation score: {:.4f}'.format(scores))
# The result is higher than 0.9, which suggest high accuracy
```

Cross-validation score: 0.9581


```
# Using GridSearch CV
parameters = [{'penalty':['l1', 'l2']],
               {'C':[1,10,100,1000]}]
grid_search = GridSearchCV(estimator = logreg,
                           param_grid = parameters,
                           scoring = 'accuracy',
                           cv = 5,
                           verbose = 0)
grid_search.fit(X_train, y_train)
```

GridSearchCV

estimator: LogisticRegression

LogisticRegression

```
# Gets the best model
print('GridSearch CV best score: {:.4f}\n\n'.format(grid_search.best_score_))
print('Parameters that give the best results: \n\n',(grid_search.best_params_))
print('\n\nEstimator that was chosen by the search: \n\n',(grid_search.best_estimator_))
```

 Cross-validation score: 0.9581

Parameters that give the best results:

```
{'penalty': 'l1'}
```

Estimator that was chosen by the search:

```
LogisticRegression(penalty='l1', random_state=0, solver='liblinear')
```

```
# Gets the GridSearch CV score on test set
print('GridSearch CV score on test set: {:.4f}\n\n'.format(grid_search.score(X_test, y_test)))
# The results are the same as with the normal model
```

GridSearch CV score on test set: 0.9581