

```
# do the same thing, but use scikitlearn randomforest classifier
```

```
!pip install scikit-learn==1.3.0 --upgrade
```

```
!pip install --upgrade xgboost
```



```
Collecting scikit-learn==1.3.0
```

```
  Downloading scikit_learn-1.3.0-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (
Requirement already satisfied: numpy>=1.17.3 in /usr/local/lib/python3.11/dist-packages (from scikit-learn)
Requirement already satisfied: scipy>=1.5.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn)
Requirement already satisfied: joblib>=1.1.1 in /usr/local/lib/python3.11/dist-packages (from scikit-learn)
Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn)
  Downloading scikit_learn-1.3.0-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (10.9 MB)
```

```
10.9/10.9 MB 35.1 MB/s eta 0:00:00
```

```
Installing collected packages: scikit-learn
```

```
  Attempting uninstall: scikit-learn
```

```
    Found existing installation: scikit-learn 1.6.1
```

```
  Uninstalling scikit-learn-1.6.1:
```

```
    Successfully uninstalled scikit-learn-1.6.1
```

```
ERROR: pip's dependency resolver does not currently take into account all the packages that are installed
imbalanced-learn 0.13.0 requires scikit-learn<2,>=1.3.2, but you have scikit-learn 1.3.0 which is incompatible.
mlxtend 0.23.4 requires scikit-learn>=1.3.1, but you have scikit-learn 1.3.0 which is incompatible.
```

```
Successfully installed scikit-learn-1.3.0
```

```
Requirement already satisfied: xgboost in /usr/local/lib/python3.11/dist-packages (2.1.4)
```

```
Requirement already satisfied: numpy in /usr/local/lib/python3.11/dist-packages (from xgboost) (1.26.4)
```

```
Requirement already satisfied: nvidia-nccl-cu12 in /usr/local/lib/python3.11/dist-packages (from xgboost) (2.19.3)
```

```
Requirement already satisfied: scipy in /usr/local/lib/python3.11/dist-packages (from xgboost) (1.13.1)
```

```
#classify with cycle features including alignment
```

```
import pandas as pd
```

```
# import xgboost as xgb
```

```
from sklearn.model_selection import train_test_split
```

```
from sklearn.ensemble import RandomForestClassifier as RFC
```

```
from sklearn.metrics import classification_report
```

```
import xgboost as xgb
```

```
from sklearn.metrics import confusion_matrix
```

```
from sklearn.metrics import roc_curve
```

```
import seaborn as sns
```

```
from matplotlib import pyplot as plt
```

```
import numpy as np
```

```
from IPython import get_ipython
```

```
from IPython.display import display
```

```
from sklearn.impute import SimpleImputer # Import SimpleImputer for imputation
```

```
import shap
```

```
shap.initjs()
```



Set up

```
df = pd.read_csv('/content/cycle_and_HMM_features_false_0-3_dataset_48days.csv')
```

```
df.head()
```

 [Show hidden output](#)

```
# LOOK AT LAUREN'S GITHUB FOR CODE

# try w xgboost
# try w subset of features
# explanatory tools to see which variables are important (SHAP values)

df = df.loc[df['pat_cat_map'].isin(['Baseline', 'PCOS'])]

df['label_01'] = df['pat_cat_map'].map({'Baseline':0, 'PCOS':1})

df = df.replace(-np.inf, np.nan)

df.columns

➡ Index(['hub_id', 'pat_cat_map', 'cycle_min', 'cycle_max', 'cycle_median',
        'cycle_mean', 'cycle_range', 'cycle_std', 'num_cycles',
        'viterbi_logprob_mean', 'viterbi_logprob_min', 'viterbi_logprob_max',
        'viterbi_logprob_std', 'viterbi_logprob_median',
        'complete_logprob_mean', 'complete_logprob_min', 'complete_logprob_max',
        'complete_logprob_std', 'complete_logprob_median', 'label_01'],
        dtype='object')


HMM_features = [ 'viterbi_logprob_mean',
                'viterbi_logprob_min', 'viterbi_logprob_max', 'viterbi_logprob_std',
                'viterbi_logprob_median', 'complete_logprob_mean',
                'complete_logprob_min', 'complete_logprob_max', 'complete_logprob_std',
                'complete_logprob_median']
cycle_features = ['cycle_min', 'cycle_max', 'cycle_median',
                 'cycle_mean', 'cycle_range', 'cycle_std']

target = 'label_01'
```


✓ All features

```
print('Performance with all features')

X_train_all, X_test_all, y_train_all, y_test_all = train_test_split(df[HMM_features+cycle_features], df[target],
                                                                    shuffle=True, random_state=51)
```

 Performance with all features

```
clf = xgb.XGBClassifier(random_state=51)
clf.fit(X_train_all, y_train_all)
y_pred_all = clf.predict(X_test_all)
y_score_all = clf.predict_proba(X_test_all)
print(confusion_matrix(y_test_all, y_pred_all, normalize='true'))
```



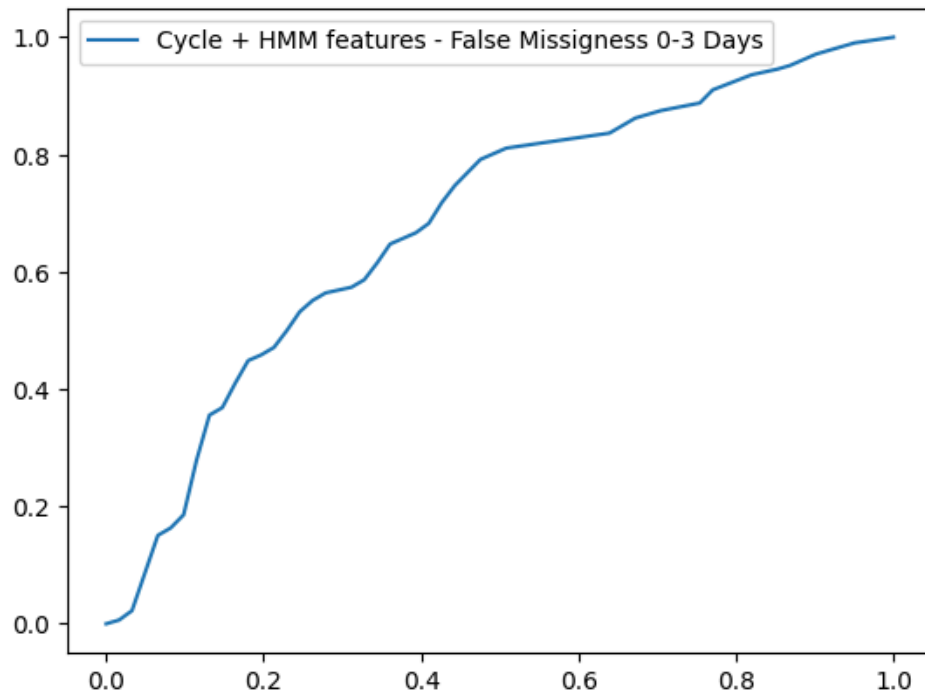
```
[[0.29508197 0.70491803]
 [0.11538462 0.88461538]]
```

```
print(classification_report(y_pred_all, y_test_all))
```

	precision	recall	f1-score	support
0	0.30	0.50	0.37	36
1	0.88	0.76	0.82	181
accuracy			0.72	217
macro avg	0.59	0.63	0.60	217
weighted avg	0.79	0.72	0.74	217

```
fpr_full, tpr_full, thresholds_full = roc_curve(y_test_all, y_score_all[:,1])#, pos_label='PCOS')
sns.lineplot(x=fpr_full, y=tpr_full, label='Cycle + HMM features - False Missigness 0-3 Days', errorbar=None)
plt.savefig('/content/drive/MyDrive/fall_research/feature distribution plots/xgb_full_features.pdf')
```

<Axes: >



```
#overall accuracy:
print((y_pred_all==y_test_all).sum()/len(y_pred_all))
```

0.7188940092165899

✓ Cycle features only

```
#PERFORMANCE WITH CYCLE FEATURES ONLY
print('Performance with cycle features only')
```

```
X_train_cycle, X_test_cycle, y_train_cycle, y_test_cycle = train_test_split(df[cycle_features], df[target],
                                                                              shuffle=True, random_state=51)
```

Performance with cycle features only

```

clf = xgb.XGBClassifier(random_state=51)
clf.fit(X_train_cycle, y_train_cycle)
y_pred_cycle = clf.predict(X_test_cycle)
y_score_cycle = clf.predict_proba(X_test_cycle)
print(confusion_matrix(y_test_cycle, y_pred_cycle, normalize='true'))

```

```

[[0.2295082 0.7704918]
 [0.1474359 0.8525641]]

```

```

print(classification_report(y_pred_cycle, y_test_cycle))

```

```

precision    recall  f1-score   support

      0       0.23      0.38      0.29         37
      1       0.85      0.74      0.79        180

 accuracy          0.68         217
 macro avg       0.54      0.56      0.54         217
 weighted avg    0.75      0.68      0.71         217

```

```

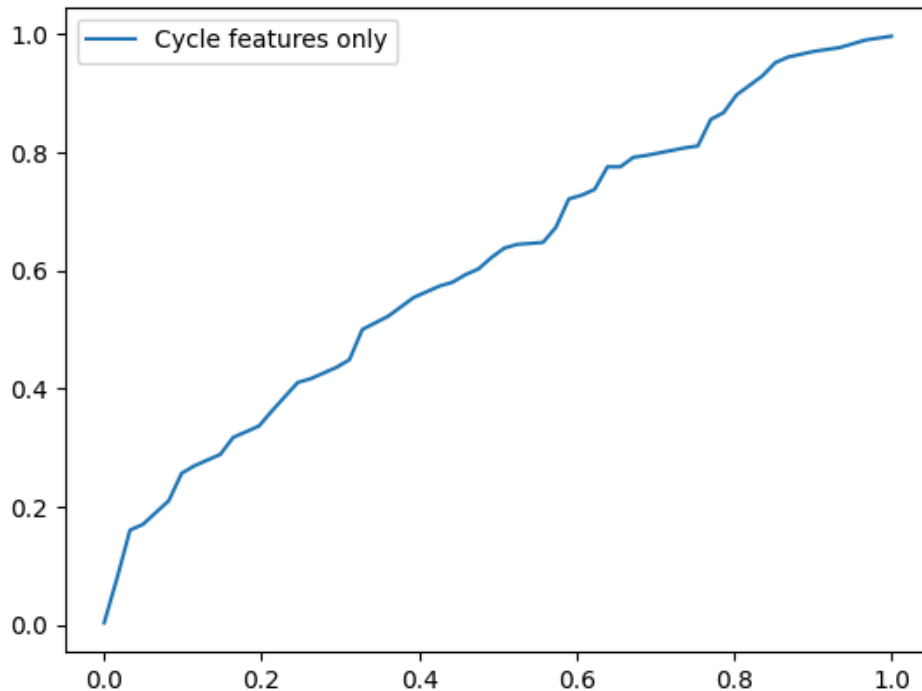
fpr_cycle, tpr_cycle, thresholds_cycle = roc_curve(y_test_cycle, y_score_cycle[:,1])#, pos_label='PCOS')
sns.lineplot(x=fpr_cycle, y=tpr_cycle, label='Cycle features only', errorbar=None)
#plt.savefig('/content/drive/MyDrive/fall_research/feature distribution plots/xgb_cycle_features_only.pdf')

```

```

<Axes: >

```



```

#overall accuracy:
print((y_pred_cycle==y_test_cycle).sum()/len(y_pred_cycle))

```

```

0.6774193548387096

```

✓ HMM Features only

```
#PERFORMANCE WITH HMM FEATURES ONLY
```

```
print('Performance with HMM features only')
```

```
X_train_hmm, X_test_hmm, y_train_hmm, y_test_hmm = train_test_split(df[HMM_features], df[target],
                                                                    shuffle=True, random_state=51)
```

↔ Performance with HMM features only

```
# Impute missing values using SimpleImputer
```

```
imputer = SimpleImputer(strategy='mean') # Replace 'mean' with other strategies if needed
```

```
X_train_hmm = imputer.fit_transform(X_train_hmm)
```

```
X_test_hmm = imputer.transform(X_test_hmm)
```

```
clf = RFC(random_state=101)
```

```
clf.fit(X_train_hmm, y_train_hmm)
```

```
y_pred_hmm = clf.predict(X_test_hmm)
```

```
y_score_hmm = clf.predict_proba(X_test_hmm)
```

```
print(confusion_matrix(y_test_hmm, y_pred_hmm, normalize='true'))
```

```
fpr_hmm, tpr_hmm, thresholds_hmm = roc_curve(y_test_hmm, y_score_hmm[:,1])#, pos_label='PCOS')
```

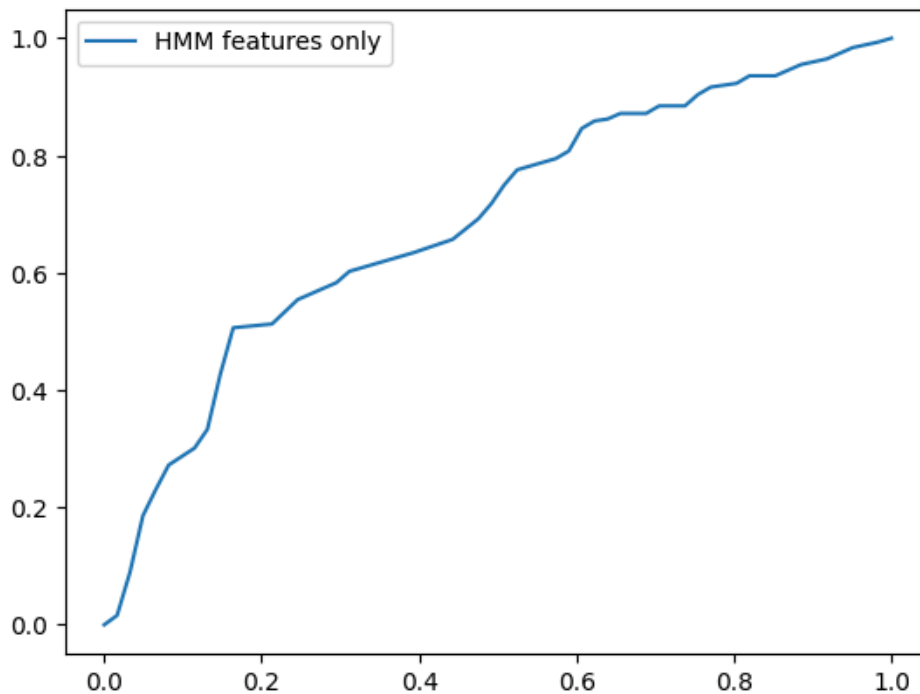
```
sns.lineplot(x=fpr_hmm, y=tpr_hmm, label='HMM features only', errorbar=None)
```

```
#plt.savefig('/content/drive/MyDrive/fall_research/feature distribution plots/xgb_hmm_features_only.pdf')
```

↔

```
[[0.2295082  0.7704918 ]
 [0.08333333 0.91666667]]
```

<Axes: >



```
print(classification_report(y_pred_cycle, y_test_cycle))
```

↔

	precision	recall	f1-score	support
0	0.23	0.38	0.29	37
1	0.85	0.74	0.79	180
accuracy			0.68	217
macro avg	0.54	0.56	0.54	217
weighted avg	0.75	0.68	0.71	217

```
#overall accuracy:
print((y_pred_cycle==y_test_cycle).sum()/len(y_pred_cycle))
#fpr_algn, tpr_algn, thresholds_algn = roc_curve(y_test, -1*X_test, pos_label='PCOS')
#sns.lineplot(x=fpr_algn, y=tpr_algn, label='HMM features only', errorbar=None)
```

↔ 0.6774193548387096

```
#make kdeplots of all features
for feature in HMM_features+cycle_features:
    sns.kdeplot(data=df, x=feature, hue='pat_cat_map', common_norm=False)
    #plt.savefig('/content/drive/MyDrive/fall_research/feature distribution plots/xgb_kdeplots_feature_disti
    plt.clf()
```

↔ <Figure size 640x480 with 0 Axes>

✓ ROC Curves

```
# put 3 ROC curves on one axis (cycle, hmm, all)
```

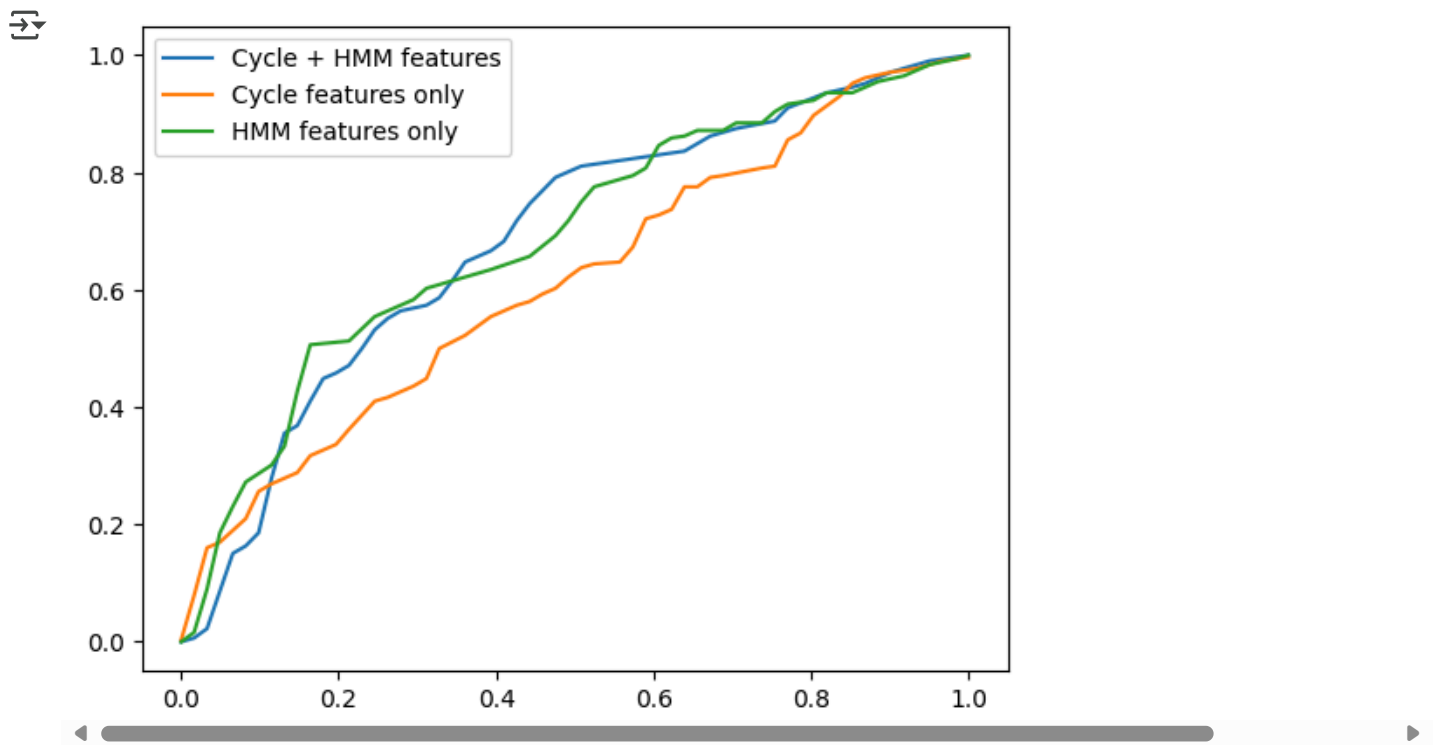
```
# # Create subplots
# fig, axes = plt.subplots(1, 3, figsize=(15, 5)) # 1 row, 3 columns
```

```
# Plot Cycle + HMM features
sns.lineplot(x=fpr_full, y=tpr_full, label='Cycle + HMM features', errorbar=None)
# axes[0].set_title("Cycle + HMM ROC Curve")
```

```
# Plot Cycle features only
sns.lineplot(x=fpr_cycle, y=tpr_cycle, label='Cycle features only', errorbar=None)
# axes[1].set_title("Cycle Only ROC Curve")
```

```
# Plot HMM features only
sns.lineplot(x=fpr_hmm, y=tpr_hmm, label='HMM features only', errorbar=None)
# axes[2].set_title("HMM Only ROC Curve")
```

```
# Adjust layout
# plt.tight_layout()
plt.show()
#plt.savefig('/content/drive/MyDrive/fall_research/feature distribution plots/xgb_roc_curves.pdf')
```



- ✓ use HMM features and take one out to see if any features are important (leave one out version)

```
HMM_features = ['viterbi_logprob_mean',
                'viterbi_logprob_min', 'viterbi_logprob_max', 'viterbi_logprob_std',
                'viterbi_logprob_median', 'complete_logprob_mean',
                'complete_logprob_min', 'complete_logprob_max', 'complete_logprob_std',
                'complete_logprob_median']
```

- ✓ without viterbi_logprob_mean

```
HMM_features = [
    'viterbi_logprob_min', 'viterbi_logprob_max', 'viterbi_logprob_std',
    'viterbi_logprob_median', 'complete_logprob_mean',
    'complete_logprob_min', 'complete_logprob_max', 'complete_logprob_std',
    'complete_logprob_median']
```

```
print('Performance with HMM features _without_viterbi_logprob_mean ')
```

```
X_train_without_viterbi_logprob_mean, X_test_without_viterbi_logprob_mean, y_train_without_viterbi_logprob_mean = train_test_split(X_train, X_test, y_train, shuffle=True, random_state=51)
```

➡ Performance with HMM features _without_viterbi_logprob_mean

```
# Impute missing values using SimpleImputer
imputer = SimpleImputer(strategy='mean') # Replace 'mean' with other strategies if needed
X_train_without_viterbi_logprob_mean = imputer.fit_transform(X_train_without_viterbi_logprob_mean)
X_test_without_viterbi_logprob_mean = imputer.transform(X_test_without_viterbi_logprob_mean)
```

```

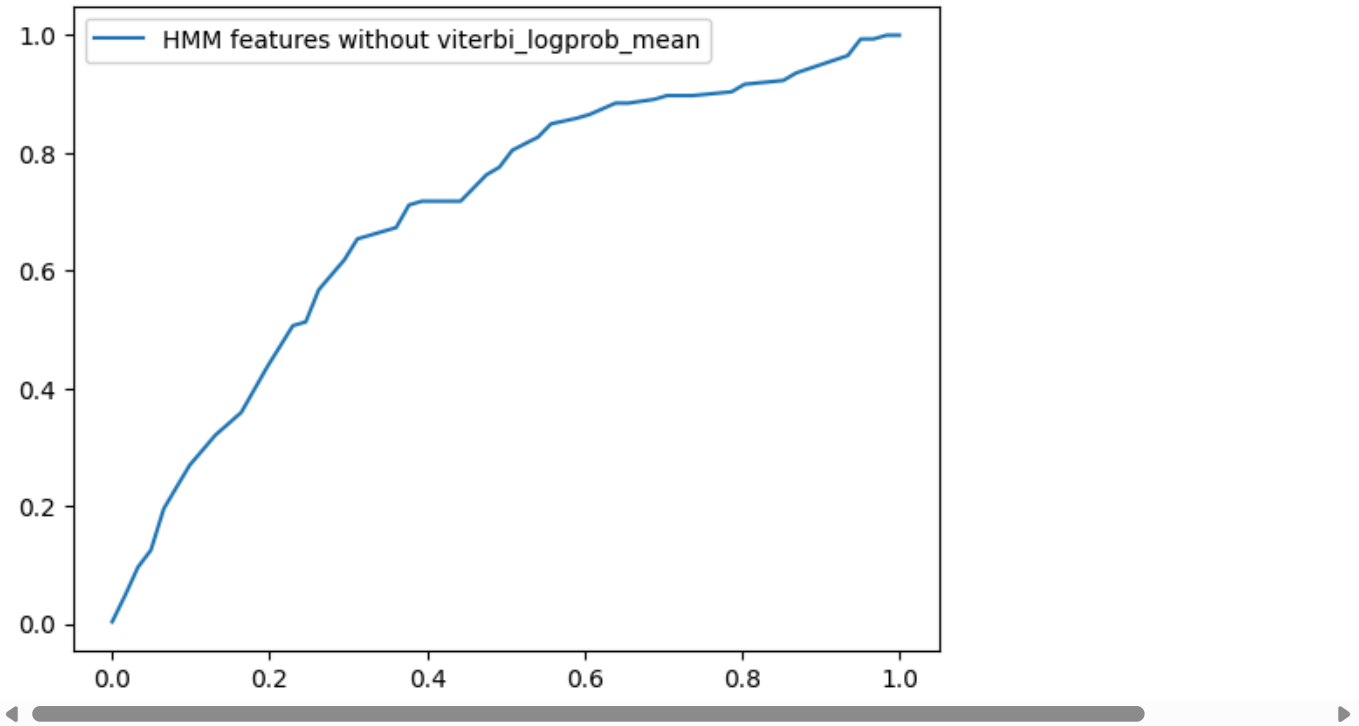
clf = RFC(random_state=101)
clf.fit(X_train_without_viterbi_logprob_mean, y_train_without_viterbi_logprob_mean)
y_pred_without_viterbi_logprob_mean = clf.predict(X_test_without_viterbi_logprob_mean)
y_score_without_viterbi_logprob_mean = clf.predict_proba(X_test_without_viterbi_logprob_mean)
print(confusion_matrix(y_test_without_viterbi_logprob_mean, y_pred_without_viterbi_logprob_mean, normalize=
fpr_without_viterbi_logprob_mean, tpr_without_viterbi_logprob_mean, thresholds_without_viterbi_logprob_mean
sns.lineplot(x=fpr_without_viterbi_logprob_mean, y=tpr_without_viterbi_logprob_mean, label='HMM features wi
#plt.savefig('/content/drive/MyDrive/fall_research/feature distribution plots/viterbi adjusted plots/xgb_wi

```

```

[[0.29508197 0.70491803]
 [0.1025641  0.8974359 ]]
<Axes: >

```



```

print(classification_report(y_pred_without_viterbi_logprob_mean, y_test_without_viterbi_logprob_mean))

```

```

precision    recall  f1-score   support

      0       0.30      0.53      0.38         34
      1       0.90      0.77      0.83        183

 accuracy          0.73         217
 macro avg       0.60      0.65      0.60         217
weighted avg       0.80      0.73      0.76         217

```

```

#overall accuracy:
print((y_pred_without_viterbi_logprob_mean==y_test_without_viterbi_logprob_mean).sum()/len(y_pred_without_
#fpr_algn, tpr_algn, thresholds_algn = roc_curve(y_test, -1*X_test, pos_label='PCOS')
#sns.lineplot(x=fpr_algn, y=tpr_algn, label='HMM features only', errorbar=None)

```

```

0.728110599078341

```

without viterbi_logprob_min


```
HMM_features = ['viterbi_logprob_mean',
                'viterbi_logprob_max', 'viterbi_logprob_std',
                'viterbi_logprob_median', 'complete_logprob_mean',
                'complete_logprob_min', 'complete_logprob_max', 'complete_logprob_std',
                'complete_logprob_median']
```

```
print('Performance with HMM features _without_viterbi_logprob_min ')
```

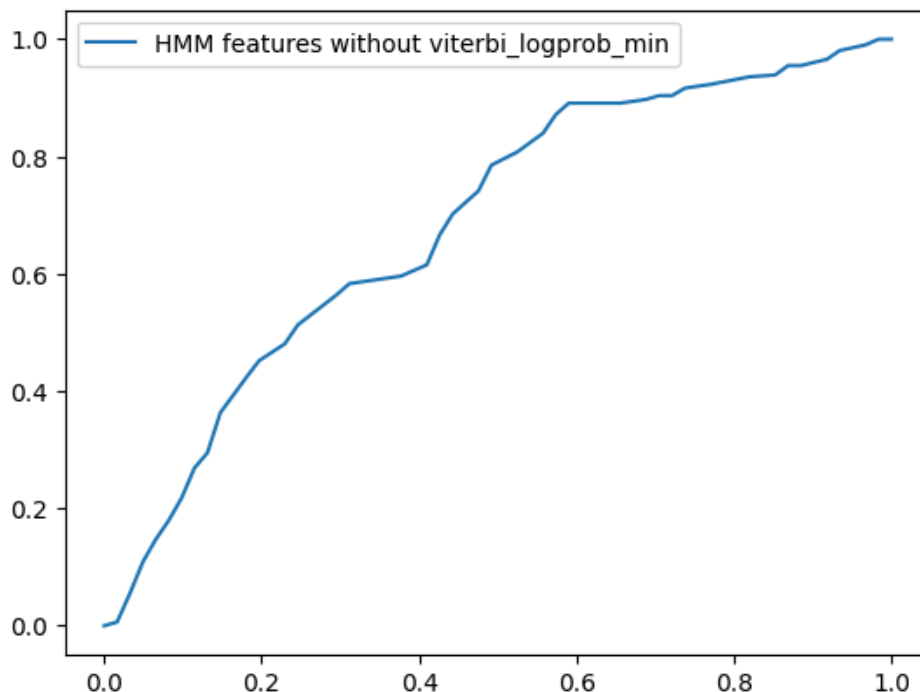
```
X_train_without_viterbi_logprob_min, X_test_without_viterbi_logprob_min, y_train_without_viterbi_logprob_min, y_test_without_viterbi_logprob_min = train_test_split(X, y, stratify=y,
                                                                    shuffle=True, random_state=51)
```

```
➡ Performance with HMM features _without_viterbi_logprob_min
```

```
# Impute missing values using SimpleImputer
imputer = SimpleImputer(strategy='mean') # Replace 'mean' with other strategies if needed
X_train_without_viterbi_logprob_min = imputer.fit_transform(X_train_without_viterbi_logprob_min)
X_test_without_viterbi_logprob_min = imputer.transform(X_test_without_viterbi_logprob_min)

clf = RFC(random_state=101)
clf.fit(X_train_without_viterbi_logprob_min, y_train_without_viterbi_logprob_min)
y_pred_without_viterbi_logprob_min = clf.predict(X_test_without_viterbi_logprob_min)
y_score_without_viterbi_logprob_min = clf.predict_proba(X_test_without_viterbi_logprob_min)
print(confusion_matrix(y_test_without_viterbi_logprob_min, y_pred_without_viterbi_logprob_min, normalize='true'))
fpr_without_viterbi_logprob_min, tpr_without_viterbi_logprob_min, thresholds_without_viterbi_logprob_min = roc_curve(y_test_without_viterbi_logprob_min, y_score_without_viterbi_logprob_min)
sns.lineplot(x=fpr_without_viterbi_logprob_min, y=tpr_without_viterbi_logprob_min, label='HMM features without viterbi_logprob_min')
plt.savefig('/content/drive/MyDrive/fall_research/feature distribution plots/viterbi adjusted plots/xgb_wi')
```

```
➡ [[0.26229508 0.73770492]
    [0.08333333 0.91666667]]
<Axes: >
```



```
print(classification_report(y_pred_without_viterbi_logprob_min, y_test_without_viterbi_logprob_min))
```

```
➡ precision    recall  f1-score   support
```

0	0.26	0.55	0.36	29
1	0.92	0.76	0.83	188
accuracy			0.73	217
macro avg	0.59	0.66	0.59	217
weighted avg	0.83	0.73	0.77	217

```
#overall accuracy:
print((y_pred_without_viterbi_logprob_min==y_test_without_viterbi_logprob_min).sum()/len(y_pred_without_viterbi_logprob_min))
#fpr_algn, tpr_algn, thresholds_algn = roc_curve(y_test, -1*X_test, pos_label='PCOS')
#sns.lineplot(x=fpr_algn, y=tpr_algn, label='HMM features only', errorbar=None)
```

→ 0.7327188940092166

✓ without viterbi_logprob_max

```
HMM_features = ['viterbi_logprob_mean',
                'viterbi_logprob_min', 'viterbi_logprob_std',
                'viterbi_logprob_median', 'complete_logprob_mean',
                'complete_logprob_min', 'complete_logprob_max', 'complete_logprob_std',
                'complete_logprob_median']

print('Performance with HMM features _without_viterbi_logprob_max ')

X_train_without_viterbi_logprob_max, X_test_without_viterbi_logprob_max, y_train_without_viterbi_logprob_max, y_test_without_viterbi_logprob_max = train_test_split(X_train, X_test, y_train, y_test,
                                                                            shuffle=True, random_state=51)

→ Performance with HMM features _without_viterbi_logprob_max

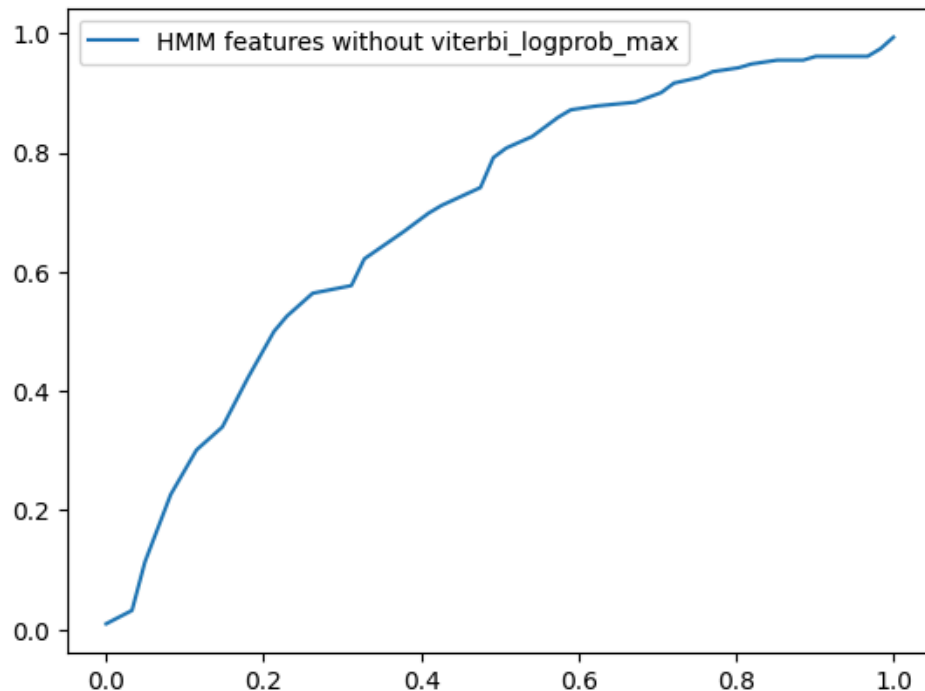
# Impute missing values using SimpleImputer
imputer = SimpleImputer(strategy='mean') # Replace 'mean' with other strategies if needed
X_train_without_viterbi_logprob_max = imputer.fit_transform(X_train_without_viterbi_logprob_max)
X_test_without_viterbi_logprob_max = imputer.transform(X_test_without_viterbi_logprob_max)

clf = RFC(random_state=101)
clf.fit(X_train_without_viterbi_logprob_max, y_train_without_viterbi_logprob_max)
y_pred_without_viterbi_logprob_max = clf.predict(X_test_without_viterbi_logprob_max)
y_score_without_viterbi_logprob_max = clf.predict_proba(X_test_without_viterbi_logprob_max)
print(confusion_matrix(y_test_without_viterbi_logprob_max, y_pred_without_viterbi_logprob_max, normalize='true'))
fpr_without_viterbi_logprob_max, tpr_without_viterbi_logprob_max, thresholds_without_viterbi_logprob_max = roc_curve(y_test_without_viterbi_logprob_max, y_score_without_viterbi_logprob_max)
sns.lineplot(x=fpr_without_viterbi_logprob_max, y=tpr_without_viterbi_logprob_max, label='HMM features without viterbi_logprob_max')
plt.savefig('/content/drive/MyDrive/fall_research/feature distribution plots/viterbi adjusted plots/xgb without viterbi_logprob_max')
```

```

[[0.24590164 0.75409836]
 [0.07692308 0.92307692]]
<Axes: >

```



```
print(classification_report(y_pred_without_viterbi_logprob_max, y_test_without_viterbi_logprob_max))
```

```

precision    recall  f1-score   support

      0       0.25      0.56      0.34         27
      1       0.92      0.76      0.83        190

   accuracy       0.73         217
  macro avg       0.58         217
weighted avg       0.84         217

```

```
#overall accuracy:
```

```

print((y_pred_without_viterbi_logprob_max==y_test_without_viterbi_logprob_max).sum()/len(y_pred_without_viterbi_logprob_max))
#fpr_algn, tpr_algn, thresholds_algn = roc_curve(y_test, -1*X_test, pos_label='PCOS')
#sns.lineplot(x=fpr_algn, y=tpr_algn, label='HMM features only', errorbar=None)

```

```
0.7327188940092166
```

without viterbi_logprob_std

```

HMM_features = ['viterbi_logprob_mean',
                'viterbi_logprob_min', 'viterbi_logprob_max',
                'viterbi_logprob_median', 'complete_logprob_mean',
                'complete_logprob_min', 'complete_logprob_max', 'complete_logprob_std',
                'complete_logprob_median']

```

```
print('Performance with HMM features _without_viterbi_logprob_std ')
```

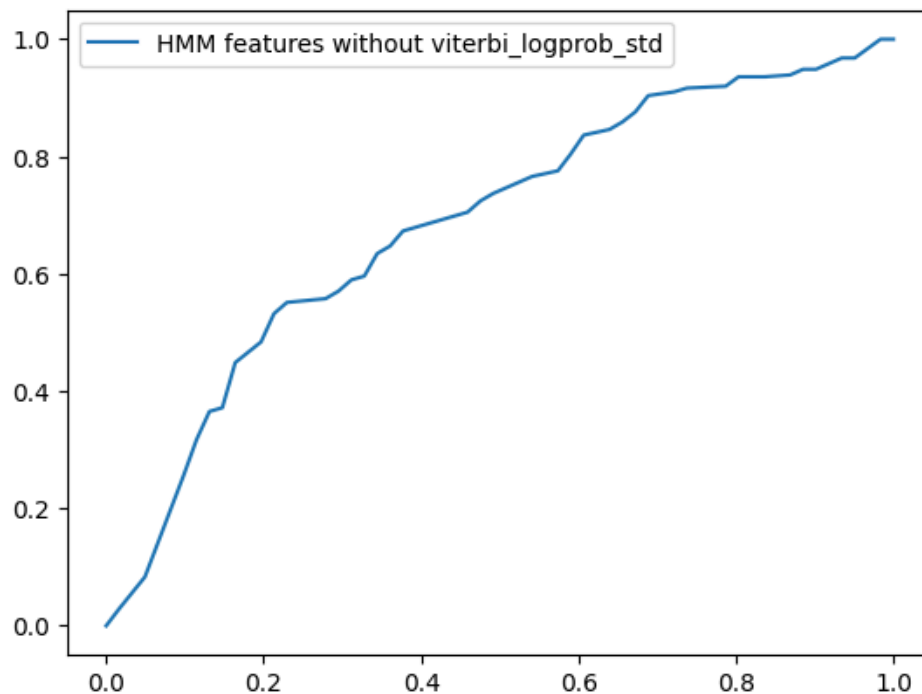
```
X_train_without_viterbi_logprob_std, X_test_without_viterbi_logprob_std, y_train_without_viterbi_logprob_std, y_test_without_viterbi_logprob_std, shuffle=True, random_state=51)
```

➡ Performance with HMM features _without_viterbi_logprob_std

```
# Impute missing values using SimpleImputer
imputer = SimpleImputer(strategy='mean') # Replace 'mean' with other strategies if needed
X_train_without_viterbi_logprob_std = imputer.fit_transform(X_train_without_viterbi_logprob_std)
X_test_without_viterbi_logprob_std = imputer.transform(X_test_without_viterbi_logprob_std)

clf = RFC(random_state=101)
clf.fit(X_train_without_viterbi_logprob_std, y_train_without_viterbi_logprob_std)
y_pred_without_viterbi_logprob_std = clf.predict(X_test_without_viterbi_logprob_std)
y_score_without_viterbi_logprob_std = clf.predict_proba(X_test_without_viterbi_logprob_std)
print(confusion_matrix(y_test_without_viterbi_logprob_std, y_pred_without_viterbi_logprob_std, normalize='true'))
fpr_without_viterbi_logprob_std, tpr_without_viterbi_logprob_std, thresholds_without_viterbi_logprob_std = roc_curve(y_test_without_viterbi_logprob_std, y_score_without_viterbi_logprob_std)
sns.lineplot(x=fpr_without_viterbi_logprob_std, y=tpr_without_viterbi_logprob_std, label='HMM features without viterbi_logprob_std')
plt.savefig('/content/drive/MyDrive/fall_research/feature distribution plots/viterbi adjusted plots/xgb_wi')
```

➡ [[0.27868852 0.72131148]
[0.08974359 0.91025641]]
<Axes: >



```
print(classification_report(y_pred_without_viterbi_logprob_std, y_test_without_viterbi_logprob_std))
```

➡

	precision	recall	f1-score	support
0	0.28	0.55	0.37	31
1	0.91	0.76	0.83	186
accuracy			0.73	217
macro avg	0.59	0.66	0.60	217
weighted avg	0.82	0.73	0.76	217

```
#overall accuracy:
print((y_pred_without_viterbi_logprob_std==y_test_without_viterbi_logprob_std).sum())/len(y_pred_without_viterbi_logprob_std)
#fpr_algn, tpr_algn, thresholds_algn = roc_curve(y_test, -1*X_test, pos_label='PCOS')
#sns.lineplot(x=fpr_algn, y=tpr_algn, label='HMM features only', errorbar=None)
```

0.7327188940092166

✓ without viterbi_logprob_median

```
HMM_features = ['viterbi_logprob_mean',
                'viterbi_logprob_min', 'viterbi_logprob_max', 'viterbi_logprob_std',
                'complete_logprob_mean',
                'complete_logprob_min', 'complete_logprob_max', 'complete_logprob_std',
                'complete_logprob_median']

print('Performance with HMM features _without_viterbi_logprob_median ')

X_train_without_viterbi_logprob_median, X_test_without_viterbi_logprob_median, y_train_without_viterbi_logprob_median = train_test_split(X_train, y_train, test_size=0.2,
                                                                    shuffle=True, random_state=51)

Performance with HMM features _without_viterbi_logprob_median

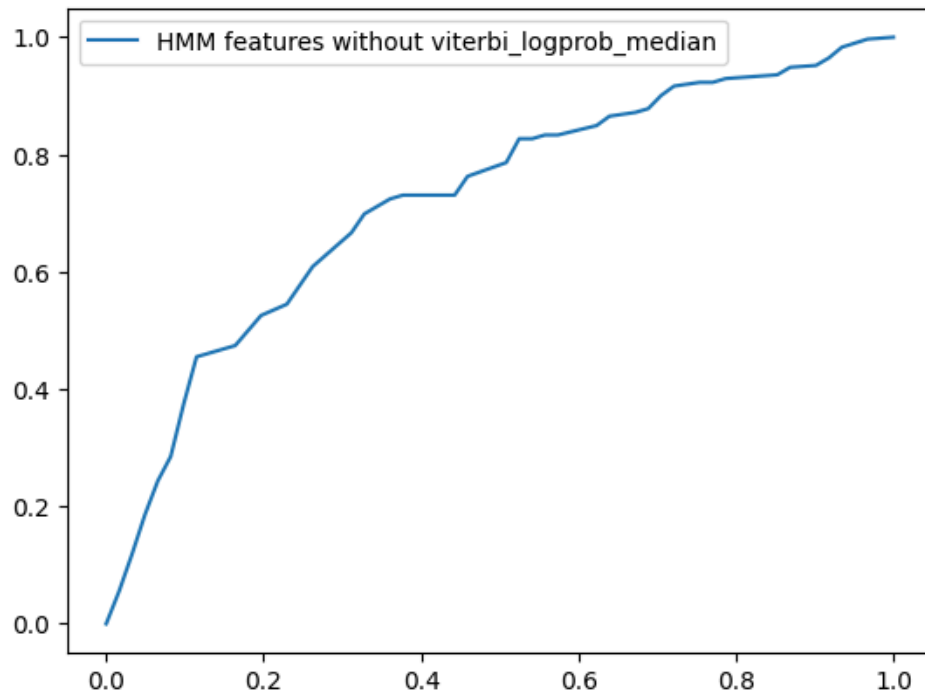
# Impute missing values using SimpleImputer
imputer = SimpleImputer(strategy='mean') # Replace 'mean' with other strategies if needed
X_train_without_viterbi_logprob_median = imputer.fit_transform(X_train_without_viterbi_logprob_median)
X_test_without_viterbi_logprob_median = imputer.transform(X_test_without_viterbi_logprob_median)

clf = RFC(random_state=101)
clf.fit(X_train_without_viterbi_logprob_median, y_train_without_viterbi_logprob_median)
y_pred_without_viterbi_logprob_median = clf.predict(X_test_without_viterbi_logprob_median)
y_score_without_viterbi_logprob_median = clf.predict_proba(X_test_without_viterbi_logprob_median)
print(confusion_matrix(y_test_without_viterbi_logprob_median, y_pred_without_viterbi_logprob_median, normalizer='true'))
fpr_without_viterbi_logprob_median, tpr_without_viterbi_logprob_median, thresholds_without_viterbi_logprob_median = roc_curve(y_test_without_viterbi_logprob_median, y_score_without_viterbi_logprob_median)
sns.lineplot(x=fpr_without_viterbi_logprob_median, y=tpr_without_viterbi_logprob_median, label='HMM features only')
plt.savefig('/content/drive/MyDrive/fall_research/feature distribution plots/viterbi adjusted plots/xgb_without_viterbi_logprob_median.png')
```

```

→ [[0.27868852 0.72131148]
   [0.07692308 0.92307692]]
<Axes: >

```



```
print(classification_report(y_pred_without_viterbi_logprob_median, y_test_without_viterbi_logprob_median))
```

```

→
              precision    recall  f1-score   support

         0       0.28        0.59        0.38         29
         1       0.92        0.77        0.84        188

 accuracy          0.74         217
 macro avg         0.60         0.68         0.61         217
 weighted avg      0.84         0.74         0.78         217

```

```
#overall accuracy:
```

```

print((y_pred_without_viterbi_logprob_median==y_test_without_viterbi_logprob_median).sum()/len(y_pred_with
#fpr_algn, tpr_algn, thresholds_algn = roc_curve(y_test, -1*X_test, pos_label='PCOS')
#sns.lineplot(x=fpr_algn, y=tpr_algn, label='HMM features only', errorbar=None)

```

```
→ 0.7419354838709677
```

✓ without complete_logprob_mean

```

HMM_features = ['viterbi_logprob_mean',
                'viterbi_logprob_min', 'viterbi_logprob_max', 'viterbi_logprob_std',
                'viterbi_logprob_median',
                'complete_logprob_min', 'complete_logprob_max', 'complete_logprob_std',
                'complete_logprob_median']

```

```
print('Performance with HMM features _without_complete_logprob_mean ')
```

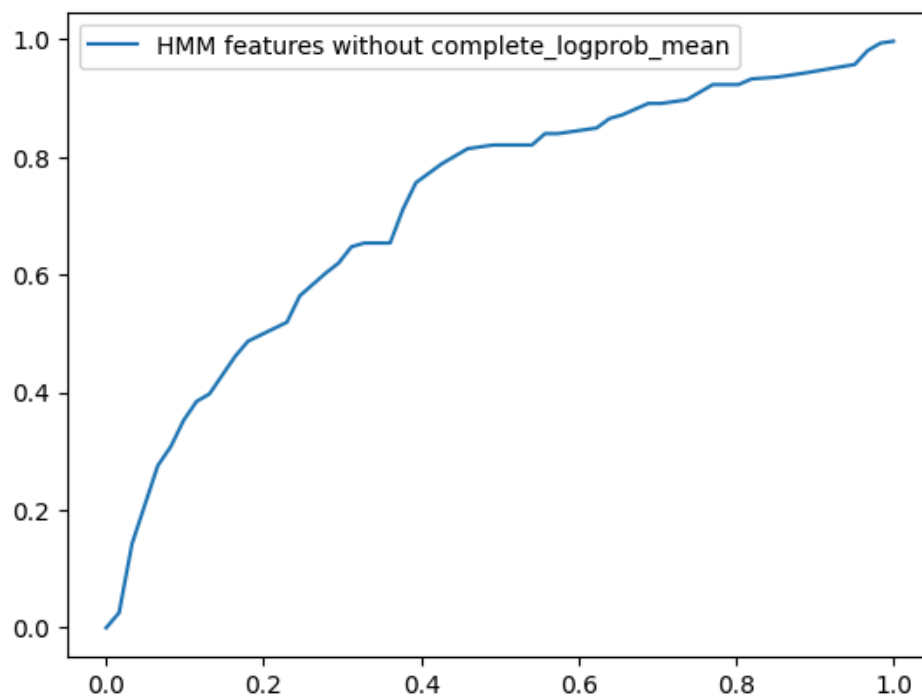
```
X_train_without_complete_logprob_mean, X_test_without_complete_logprob_mean, y_train_without_complete_logprob_mean, y_test_without_complete_logprob_mean)
shuffle=True, random_state=51)
```

➡ Performance with HMM features _without_complete_logprob_mean

```
# Impute missing values using SimpleImputer
imputer = SimpleImputer(strategy='mean') # Replace 'mean' with other strategies if needed
X_train_without_complete_logprob_mean = imputer.fit_transform(X_train_without_complete_logprob_mean)
X_test_without_complete_logprob_mean = imputer.transform(X_test_without_complete_logprob_mean)

clf = RFC(random_state=101)
clf.fit(X_train_without_complete_logprob_mean, y_train_without_complete_logprob_mean)
y_pred_without_complete_logprob_mean = clf.predict(X_test_without_complete_logprob_mean)
y_score_without_complete_logprob_mean = clf.predict_proba(X_test_without_complete_logprob_mean)
print(confusion_matrix(y_test_without_complete_logprob_mean, y_pred_without_complete_logprob_mean, normalize=True))
fpr_without_complete_logprob_mean, tpr_without_complete_logprob_mean, thresholds_without_complete_logprob_mean = roc_curve(y_test_without_complete_logprob_mean, y_score_without_complete_logprob_mean)
sns.lineplot(x=fpr_without_complete_logprob_mean, y=tpr_without_complete_logprob_mean, label='HMM features without complete_logprob_mean')
plt.savefig('/content/drive/MyDrive/fall_research/feature distribution plots/viterbi adjusted plots/xgb_wi')
```

➡
[[0.24590164 0.75409836]
[0.09615385 0.90384615]]
<Axes: >



```
print(classification_report(y_pred_without_complete_logprob_mean, y_test_without_complete_logprob_mean))
```

➡

	precision	recall	f1-score	support
0	0.25	0.50	0.33	30
1	0.90	0.75	0.82	187
accuracy			0.72	217
macro avg	0.57	0.63	0.58	217
weighted avg	0.81	0.72	0.75	217

```
#overall accuracy:
print((y_pred_without_complete_logprob_mean==y_test_without_complete_logprob_mean).sum())/len(y_pred_without
#fpr_algn, tpr_algn, thresholds_algn = roc_curve(y_test, -1*X_test, pos_label='PCOS')
#sns.lineplot(x=fpr_algn, y=tpr_algn, label='HMM features only', errorbar=None)

➡ 0.7188940092165899
```

✓ without complete_logprob_min

```
HMM_features = ['viterbi_logprob_mean',
                'viterbi_logprob_min', 'viterbi_logprob_max', 'viterbi_logprob_std',
                'viterbi_logprob_median', 'complete_logprob_mean',
                'complete_logprob_max', 'complete_logprob_std',
                'complete_logprob_median']

print('Performance with HMM features _without_complete_logprob_min ')

X_train_without_complete_logprob_min, X_test_without_complete_logprob_min, y_train_without_complete_logpro
shuffle=True, random_state=51)

➡ Performance with HMM features _without_complete_logprob_min

# Impute missing values using SimpleImputer
imputer = SimpleImputer(strategy='mean') # Replace 'mean' with other strategies if needed
X_train_without_complete_logprob_min = imputer.fit_transform(X_train_without_complete_logprob_min)
X_test_without_complete_logprob_min = imputer.transform(X_test_without_complete_logprob_min)

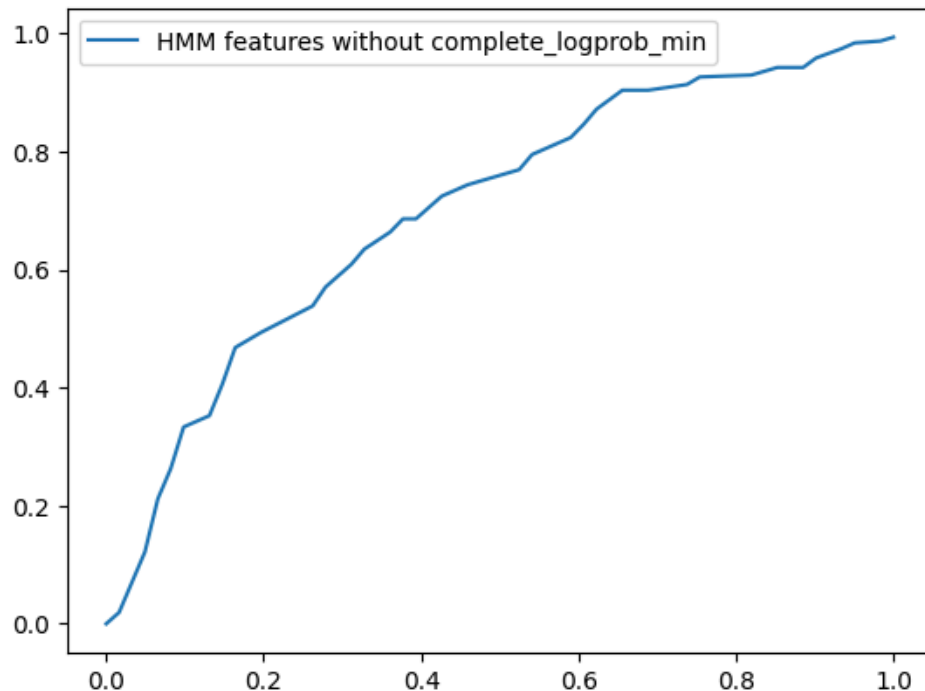
clf = RFC(random_state=101)
clf.fit(X_train_without_complete_logprob_min, y_train_without_complete_logprob_min)
y_pred_without_complete_logprob_min = clf.predict(X_test_without_complete_logprob_min)
y_score_without_complete_logprob_min = clf.predict_proba(X_test_without_complete_logprob_min)
print(confusion_matrix(y_test_without_complete_logprob_min, y_pred_without_complete_logprob_min, normalize=
fpr_without_complete_logprob_min, tpr_without_complete_logprob_min, thresholds_without_complete_logprob_min
sns.lineplot(x=fpr_without_complete_logprob_min, y=tpr_without_complete_logprob_min, label='HMM features wi
#plt.savefig('/content/drive/MyDrive/fall_research/feature distribution plots/viterbi adjusted plots/xgb_wi
```



```

→ [[0.26229508 0.73770492]
   [0.08333333 0.91666667]]
<Axes: >

```



```
print(classification_report(y_pred_without_complete_logprob_min, y_test_without_complete_logprob_min))
```

```

→
              precision    recall  f1-score   support

     0       0.26         0.55         0.36         29
     1       0.92         0.76         0.83        188

 accuracy          0.73         217
 macro avg         0.59         0.66         0.59         217
 weighted avg      0.83         0.73         0.77         217

```

```
#overall accuracy:
```

```

print((y_pred_without_complete_logprob_min==y_test_without_complete_logprob_min).sum()/len(y_pred_without_
#fpr_algn, tpr_algn, thresholds_algn = roc_curve(y_test, -1*X_test, pos_label='PCOS')
#sns.lineplot(x=fpr_algn, y=tpr_algn, label='HMM features only', errorbar=None)

```

```
→ 0.7327188940092166
```

✓ without complete_logprob_max

```

HMM_features = ['viterbi_logprob_mean',
                'viterbi_logprob_min', 'viterbi_logprob_max', 'viterbi_logprob_std',
                'viterbi_logprob_median', 'complete_logprob_mean',
                'complete_logprob_min', 'complete_logprob_std',
                'complete_logprob_median']

```

```
print('Performance with HMM features _without_complete_logprob_max ')
```

```
X_train_without_complete_logprob_max, X_test_without_complete_logprob_max, y_train_without_complete_logprob_max, y_test_without_complete_logprob_max)
shuffle=True, random_state=51)
```

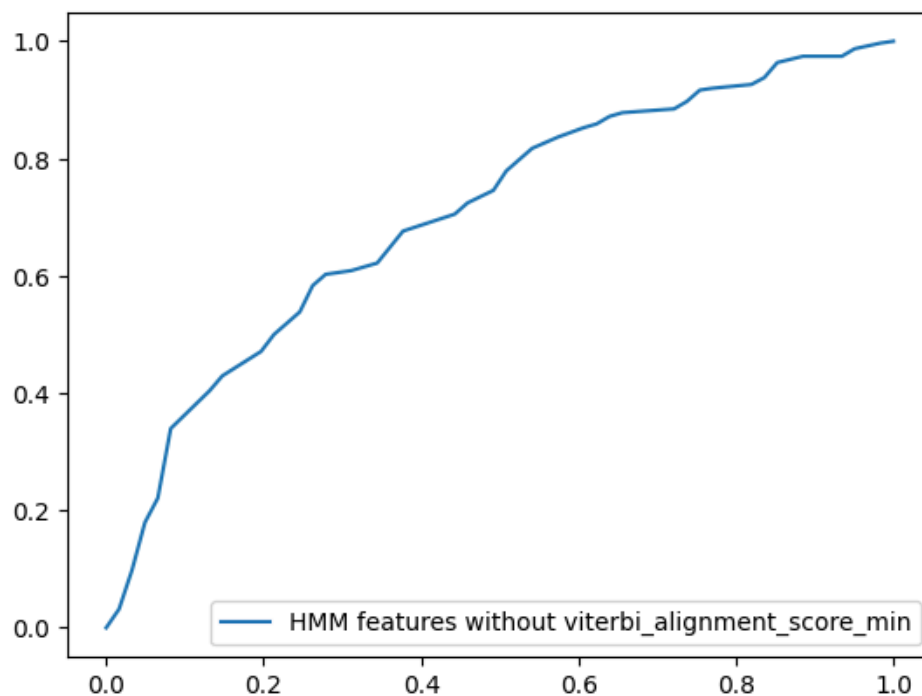
➡ Performance with HMM features _without_complete_logprob_max

```
# Impute missing values using SimpleImputer
imputer = SimpleImputer(strategy='mean') # Replace 'mean' with other strategies if needed
X_train_without_complete_logprob_max = imputer.fit_transform(X_train_without_complete_logprob_max)
X_test_without_complete_logprob_max = imputer.transform(X_test_without_complete_logprob_max)

clf = RFC(random_state=101)
clf.fit(X_train_without_complete_logprob_max, y_train_without_complete_logprob_max)
y_pred_without_complete_logprob_max = clf.predict(X_test_without_complete_logprob_max)
y_score_without_complete_logprob_max = clf.predict_proba(X_test_without_complete_logprob_max)
print(confusion_matrix(y_test_without_complete_logprob_max, y_pred_without_complete_logprob_max, normalize=
fpr_without_complete_logprob_max, tpr_without_complete_logprob_max, thresholds_without_complete_logprob_max)
sns.lineplot(x=fpr_without_complete_logprob_max, y=tpr_without_complete_logprob_max, label='HMM features wi
#plt.savefig('/content/drive/MyDrive/fall_research/feature distribution plots/viterbi adjusted plots/xgb_wi
```

➡

```
[[0.24590164 0.75409836]
 [0.08333333 0.91666667]]
<Axes: >
```



```
print(classification_report(y_pred_without_complete_logprob_max, y_test_without_complete_logprob_max))
```

➡

	precision	recall	f1-score	support
0	0.25	0.54	0.34	28
1	0.92	0.76	0.83	189
accuracy			0.73	217
macro avg	0.58	0.65	0.58	217
weighted avg	0.83	0.73	0.77	217

```
#overall accuracy:
print((y_pred_without_complete_logprob_max==y_test_without_complete_logprob_max).sum()/len(y_pred_without_
#fpr_algn, tpr_algn, thresholds_algn = roc_curve(y_test, -1*X_test, pos_label='PCOS')
#sns.lineplot(x=fpr_algn, y=tpr_algn, label='HMM features only', errorbar=None)

➡ 0.728110599078341
```

✓ without complete_logprob_std

```
HMM_features = ['viterbi_logprob_mean',
                'viterbi_logprob_min', 'viterbi_logprob_max', 'viterbi_logprob_std',
                'viterbi_logprob_median', 'complete_logprob_mean',
                'complete_logprob_min', 'complete_logprob_max',
                'complete_logprob_median']

print('Performance with HMM features _without_complete_logprob_std ')

X_train_without_complete_logprob_std, X_test_without_complete_logprob_std, y_train_without_complete_logpro
shuffle=True, random_state=51)

➡ Performance with HMM features _without_complete_logprob_std

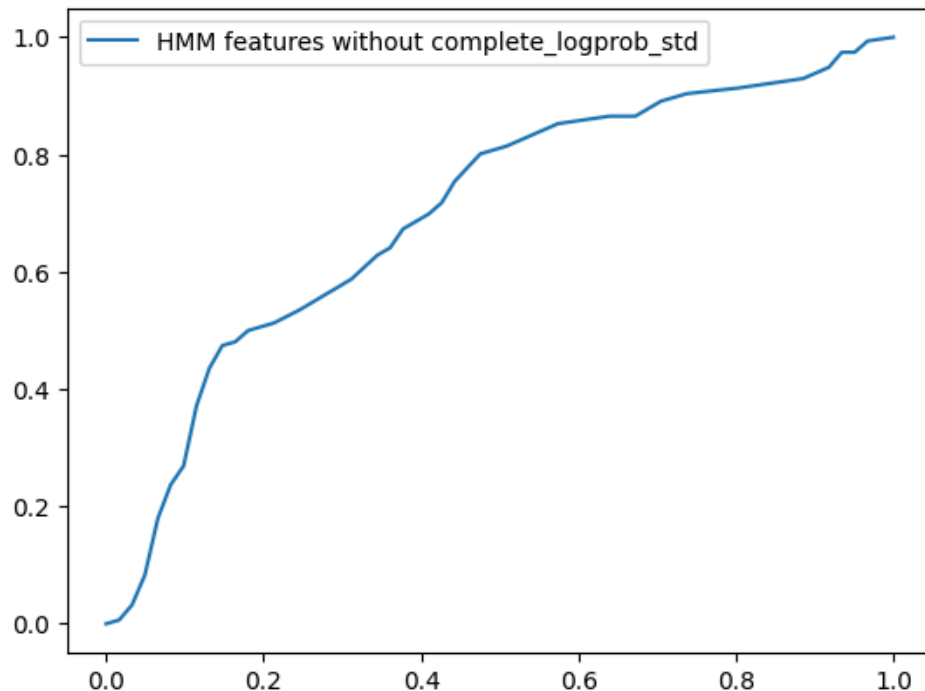
# Impute missing values using SimpleImputer
imputer = SimpleImputer(strategy='mean') # Replace 'mean' with other strategies if needed
X_train_without_complete_logprob_std = imputer.fit_transform(X_train_without_complete_logprob_std)
X_test_without_complete_logprob_std = imputer.transform(X_test_without_complete_logprob_std)

clf = RFC(random_state=101)
clf.fit(X_train_without_complete_logprob_std, y_train_without_complete_logprob_std)
y_pred_without_complete_logprob_std = clf.predict(X_test_without_complete_logprob_std)
y_score_without_complete_logprob_std = clf.predict_proba(X_test_without_complete_logprob_std)
print(confusion_matrix(y_test_without_complete_logprob_std, y_pred_without_complete_logprob_std, normalize=
fpr_without_complete_logprob_std, tpr_without_complete_logprob_std, thresholds_without_complete_logprob_std
sns.lineplot(x=fpr_without_complete_logprob_std, y=tpr_without_complete_logprob_std, label='HMM features wi
#plt.savefig('/content/drive/MyDrive/fall_research/feature distribution plots/viterbi adjusted plots/xgb_wi
```

```

[[0.26229508 0.73770492]
 [0.09615385 0.90384615]]
<Axes: >

```



```
print(classification_report(y_pred_without_complete_logprob_std, y_test_without_complete_logprob_std))
```

```

precision    recall  f1-score   support

      0       0.26      0.52      0.35         31
      1       0.90      0.76      0.82        186

 accuracy          0.72         217
 macro avg          0.58         217
weighted avg          0.81         217

```

```

#overall accuracy:
print((y_pred_without_complete_logprob_std==y_test_without_complete_logprob_std).sum()/len(y_pred_without_
#fpr_algn, tpr_algn, thresholds_algn = roc_curve(y_test, -1*X_test, pos_label='PCOS')
#sns.lineplot(x=fpr_algn, y=tpr_algn, label='HMM features only', errorbar=None)

```

```
0.7235023041474654
```

without complete_logprob_median

```

HMM_features = ['viterbi_logprob_mean',
                'viterbi_logprob_min', 'viterbi_logprob_max', 'viterbi_logprob_std',
                'viterbi_logprob_median', 'complete_logprob_mean',
                'complete_logprob_min', 'complete_logprob_max', 'complete_logprob_std']

```

```
print('Performance with HMM features _without_viterbi_logprob_median ')
```

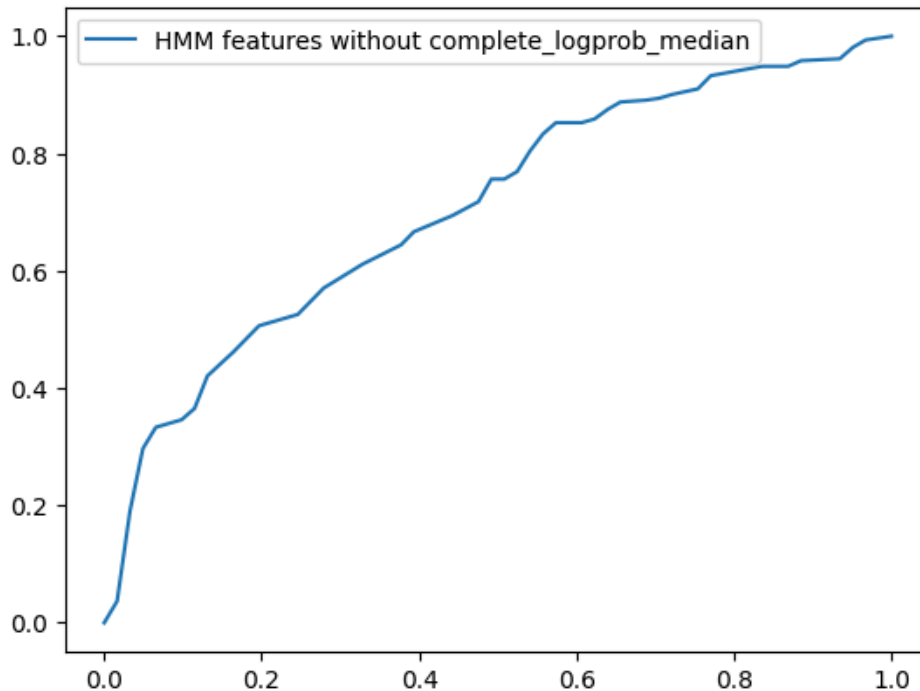
```
X_train_without_viterbi_logprob_median, X_test_without_viterbi_logprob_median, y_train_without_viterbi_logprob_median, y_test_without_viterbi_logprob_median = train_test_split(X_train_without_viterbi_logprob_median, y_train_without_viterbi_logprob_median, X_test_without_viterbi_logprob_median, y_test_without_viterbi_logprob_median, shuffle=True, random_state=51)
```

```
→ Performance with HMM features _without_viterbi_logprob_median
```

```
# Impute missing values using SimpleImputer
imputer = SimpleImputer(strategy='mean') # Replace 'mean' with other strategies if needed
X_train_without_viterbi_logprob_median = imputer.fit_transform(X_train_without_viterbi_logprob_median)
X_test_without_viterbi_logprob_median = imputer.transform(X_test_without_viterbi_logprob_median)

clf = RFC(random_state=101)
clf.fit(X_train_without_viterbi_logprob_median, y_train_without_viterbi_logprob_median)
y_pred_without_viterbi_logprob_median = clf.predict(X_test_without_viterbi_logprob_median)
y_score_without_viterbi_logprob_median = clf.predict_proba(X_test_without_viterbi_logprob_median)
print(confusion_matrix(y_test_without_viterbi_logprob_median, y_pred_without_viterbi_logprob_median, normalizer=y_score_without_viterbi_logprob_median))
fpr_without_viterbi_logprob_median, tpr_without_viterbi_logprob_median, thresholds_without_viterbi_logprob_median = roc_curve(y_test_without_viterbi_logprob_median, y_score_without_viterbi_logprob_median)
sns.lineplot(x=fpr_without_viterbi_logprob_median, y=tpr_without_viterbi_logprob_median, label='HMM feature: _without_viterbi_logprob_median')
plt.savefig('/content/drive/MyDrive/fall_research/feature_distribution_plots/viterbi_adjusted_plots/xgb_wi')
```

```
→ [[0.27868852 0.72131148]
     [0.1025641  0.8974359 ]]
<Axes: >
```



```
print(classification_report(y_pred_without_viterbi_logprob_median, y_test_without_viterbi_logprob_median))
```

```
→
```

	precision	recall	f1-score	support
0	0.28	0.52	0.36	33
1	0.90	0.76	0.82	184
accuracy			0.72	217
macro avg	0.59	0.64	0.59	217
weighted avg	0.80	0.72	0.75	217

```
#overall accuracy:
print((y_pred_without_viterbi_logprob_median==y_test_without_viterbi_logprob_median).sum()/len(y_pred_with
#fpr_algn, tpr_algn, thresholds_algn = roc_curve(y_test, -1*X_test, pos_label='PCOS')
#sns.lineplot(x=fpr_algn, y=tpr_algn, label='HMM features only', errorbar=None)
```

→ 0.7235023041474654

```
HMM_features = ['viterbi_logprob_mean',
                'viterbi_logprob_min', 'viterbi_logprob_max', 'viterbi_logprob_std',
                'viterbi_logprob_median', 'complete_logprob_mean',
                'complete_logprob_min', 'complete_logprob_max', 'complete_logprob_std']
```

```
print('Performance with HMM features _without_viterbi_alignment ')
```

```
X_train_without_viterbi_alignment, X_test_without_viterbi_alignment, y_train_without_viterbi_alignment, y_
shuffle=True, random_state=51)
```

→ Performance with HMM features _without_viterbi_alignment

```
# Impute missing values using SimpleImputer
imputer = SimpleImputer(strategy='mean') # Replace 'mean' with other strategies if needed
X_train_without_viterbi_alignment = imputer.fit_transform(X_train_without_viterbi_alignment)
X_test_without_viterbi_alignment = imputer.transform(X_test_without_viterbi_alignment)
```

```
clf = RFC(random_state=101)
clf.fit(X_train_without_viterbi_alignment, y_train_without_viterbi_alignment)
y_pred_without_viterbi_alignment = clf.predict(X_test_without_viterbi_alignment)
y_score_without_viterbi_alignment = clf.predict_proba(X_test_without_viterbi_alignment)
print(confusion_matrix(y_test_without_viterbi_alignment, y_pred_without_viterbi_alignment, normalize='true'
fpr_without_viterbi_alignment, tpr_without_viterbi_alignment, thresholds_without_viterbi_alignment = roc_cui
sns.lineplot(x=fpr_without_viterbi_alignment, y=tpr_without_viterbi_alignment, label='HMM features without ;
#plt.savefig('/content/drive/MyDrive/fall_research/feature distribution plots/viterbi adjusted plots/xgb_wi
```

→ [[0.27868852 0.72131148]
[0.1025641 0.8974359]]
<Axes: >

