

# Business analytics homework assignment

Student ID: 710045253

**Note: The data used in this assignment: mailing\_test & mailing\_train (not 1k & 4k)**

## 1. Fit a Random Forest model

```
In [1]: 1 import numpy as np
2 import matplotlib.pyplot as plt
3 import statsmodels.api as sm
4 import pandas as pd
5 from sklearn import datasets
6 import csv
7 import pandas as pd
```

Our outcome variable is class. Provide a quick look at the distribution of class for the training data and test data.

```
In [2]: 1 #Reading .csv file from mailing_train and mailing_test, removing the first
2 train_df = pd.read_csv('mailing_train.csv')
3 del train_df[train_df.columns[0]]
4 test_df = pd.read_csv('mailing_test.csv')
5 del test_df[test_df.columns[0]]
6
7 #Counting the number of donations(class) in mailing train and mailing test
8 train_class_counts = train_df['class'].value_counts()
9 print(train_class_counts)
10
11 test_class_counts = test_df['class'].value_counts()
12 print(test_class_counts)
```

```
1    10000
0    10000
Name: class, dtype: int64
0     4738
1      262
Name: class, dtype: int64
```

**Comparing the distribution of the outcome variable in training and test, do they look balanced?**

The training set is perfectly balanced(50/50). The testing set looks imbalanced, but it won't impact the model training process.

**Fit a random forest model to predict class.**

```

In [4]: 1 #define X_train, X_test, y_train, y_test
2 X_train = train_df.drop('class', axis=1)
3 y_train = train_df['class']
4
5 X_test = test_df.drop('class', axis=1)
6 y_test = test_df['class']
7
8 #X_train.head()

In [5]: 1 from sklearn.ensemble import RandomForestClassifier
2
3 #instantiate the model
4 rf = RandomForestClassifier(random_state=0, oob_score = True, n_estimators
5
6 #fit the model with data
7 rf.fit(X_train, y_train)
8
9 #evaluate model
10 def evaluate_model(model, X_test, y_test):
11     from sklearn import metrics
12     import seaborn as sns
13
14 rf_eval = evaluate_model(rf,X_train,y_train)

```

Examine the results.

```

In [6]: 1 print('OOB Score', rf.oob_score_)
2 print('OOB error', 1 - rf.oob_score_)

```

OOB Score 0.756

OOB error 0.244

What is the out-of-box error rate?

OOB estimate of error rate is 24.4%, showing an effective prediction.

**2. Compute and Compare Predictive Performance**

Use the *confusionMatrix* function to compute several metrics of predictive performance.

```

In [7]: 1 def evaluate_model(model, X_test, y_test):
2         from sklearn import metrics
3         import seaborn as sns
4         #making predictions
5         y_pred = model.predict(X_test)
6
7         #Computing evaluation metrics
8         acc = metrics.accuracy_score(y_test, y_pred)
9         prec = metrics.precision_score(y_test, y_pred)
10        rec = metrics.recall_score(y_test, y_pred)
11        f1 = metrics.f1_score(y_test, y_pred)
12        kappa = metrics.cohen_kappa_score(y_test, y_pred)
13
14        #Calculating Area Under Curve (AUC)
15        y_pred_proba = model.predict_proba(X_test)[:,-1]
16        fpr, tpr, _ = metrics.roc_curve(y_test, y_pred_proba)
17        auc = metrics.roc_auc_score(y_test, y_pred_proba)
18
19        #Display Confusion Matrix
20        cm = metrics.confusion_matrix(y_test, y_pred)
21
22        #Calculating Specificity (TNR) = TN/(TN+FP)
23        FP = cm[0][1]
24        TN = cm[0][0]
25        TNR = TN/(TN+FP)
26
27        return {'acc': acc, 'prec': prec, 'rec': rec, 'f1': f1, 'kappa': kappa,
28                'cm': cm, 'TNR':TNR}
29
30 rf_eval = evaluate_model(rf,X_train,y_train)
31 rf_eval_test = evaluate_model(rf, X_test, y_test)
32
33 #Print Metrics
34 print('Accuracy:', rf_eval['acc'])
35 print('Precision:', rf_eval['prec'])
36 print('Recall:', rf_eval['rec'])
37 print('F1 Score:', rf_eval['f1'])
38 print('Cohens Kappa Score:', rf_eval['kappa'])
39 print('Area Under Curve:', rf_eval['auc'])
40 print('Confusion Matrix:\n', rf_eval['cm'])
41 print('Specificity:', rf_eval['TNR'])

```

```

Accuracy: 0.98125
Precision: 0.9850821489769177
Recall: 0.9773
F1 Score: 0.9811756437929823
Cohens Kappa Score: 0.9625
Area Under Curve: 0.998277995
Confusion Matrix:
[[9852 148]
 [ 227 9773]]
Specificity: 0.9852

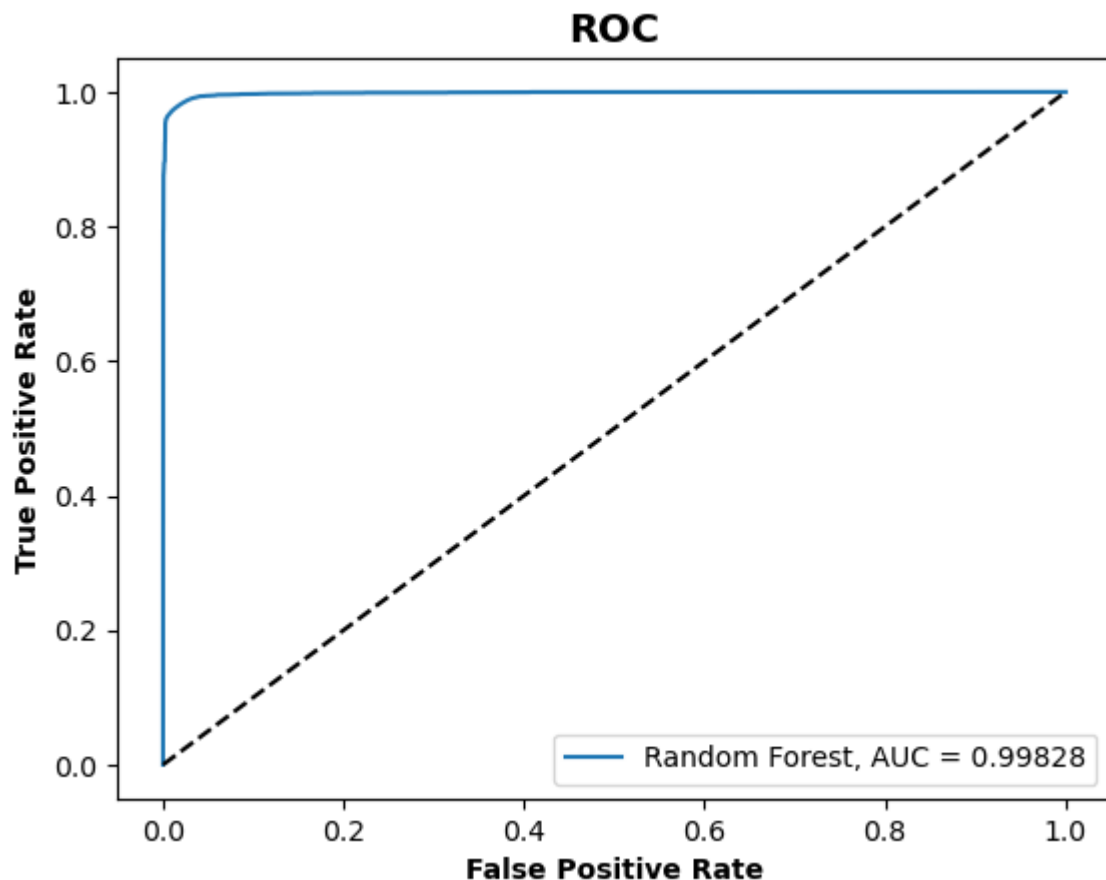
```

How would you describe the performance of this model?

The accuracy is 98.1%, which is pretty good. The model is more specific(98.52%) than sensitive(97.73%). The model is slightly better at predicting people not giving a donation than people giving donations.

The confusion matrix assumed a 0.50 cutoff for prediction. Now create a ROC plot and compute the AUC for the training set

```
In [8]: 1 #evaluate training set
2 rf_eval = evaluate_model(rf,X_train,y_train)
3 # plot roc curve
4 plt.plot(rf_eval['fpr'], rf_eval['tpr'],
5          label='Random Forest, AUC = {:.5f}'.format(rf_eval['auc']))
6
7 # Configure x and y axis
8 plt.xlabel('False Positive Rate', fontweight='bold')
9 plt.ylabel('True Positive Rate', fontweight='bold')
10 plt.plot([0, 1], [0, 1], 'k--')
11
12 # Create legend & title
13 plt.title('ROC', fontsize=14, fontweight='bold')
14 plt.legend(loc=4)
15
16 plt.show()
```

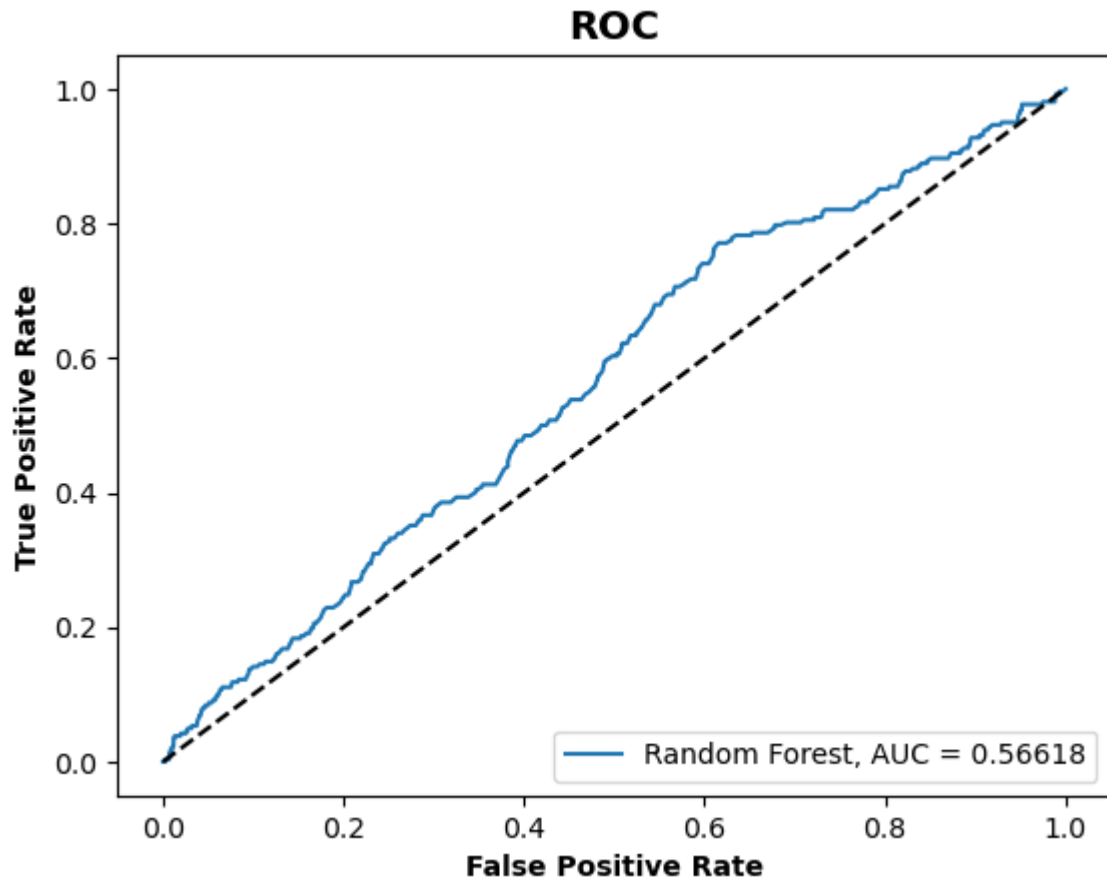


Now create a ROC chart for the test set and compute the test AUC.

```

In [9]: 1 # plot roc curve
2 plt.plot(rf_eval_test['fpr'], rf_eval_test['tpr'],
3          label='Random Forest, AUC = {:.5f}'.format(rf_eval_test['auc']))
4
5 # Configure x and y axis
6 plt.xlabel('False Positive Rate', fontweight='bold')
7 plt.ylabel('True Positive Rate', fontweight='bold')
8 plt.plot([0, 1], [0, 1], 'k--')
9
10 # Create Legend & title
11 plt.title('ROC', fontsize=14, fontweight='bold')
12 plt.legend(loc=4)
13
14 plt.show()
15
16 #Print Metrics
17 print('Accuracy:', rf_eval_test['acc'])
18 print('Precision:', rf_eval_test['prec'])
19 print('Recall:', rf_eval_test['rec'])
20 print('F1 Score:', rf_eval_test['f1'])
21 print('Cohens Kappa Score:', rf_eval_test['kappa'])
22 print('Area Under Curve:', rf_eval_test['auc'])
23 print('Confusion Matrix:\n', rf_eval_test['cm'])
24 print('Specificity:', rf_eval_test['TNR'])

```



```
Accuracy: 0.706
Precision: 0.06609195402298851
Recall: 0.3511450381679389
F1 Score: 0.11124546553808949
Cohens Kappa Score: 0.025275469087838065
Area Under Curve: 0.5661760204163834
Confusion Matrix:
[[3438 1300]
 [ 170   92]]
Specificity: 0.7256226255804137
```

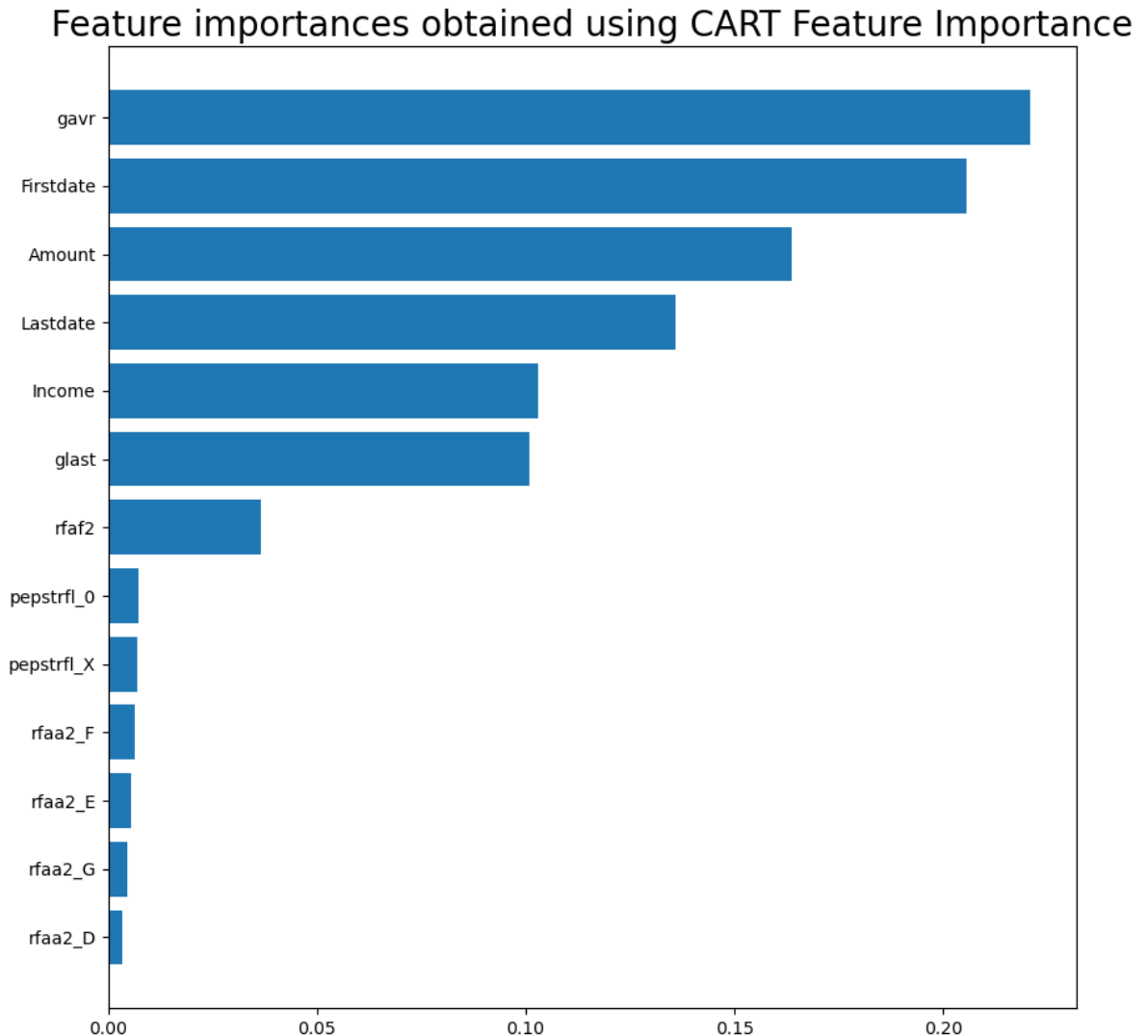
Is the model underfit, overfit, or correctly fit to the data?

The model is overfit to the data. The test-AUC is 0.4321 units lower than the train-AUC, indicating that the model does not correctly fit to the data.

### 3. Examine the model

Examine the variable importance of the model

```
In [10]: 1 sorted_idx = rf.feature_importances_.argsort()
2 plt.figure(figsize=(10,10))
3 plt.barh(X_train.columns[sorted_idx],rf.feature_importances_[sorted_idx])
4 plt.title('Feature importances obtained using CART Feature Importance', si
5 plt.show()
```



Makde some individual predictions of the model. Choose one case from the data and see what the model predicts for that one person.

```
In [11]: 1 d = X_train.loc[15560] #grabbing the 15559th data
2 d_resaped = d.values.reshape(1, -1) #changing to 2-D
3 rf.predict_proba(d_resaped)
```

D:\Programming\Anaconda\lib\site-packages\sklearn\base.py:450: UserWarning: X does not have valid feature names, but RandomForestClassifier was fitted with feature names

```
warnings.warn(
```

```
Out[11]: array([[0.79, 0.21]])
```

Using the most important variable from the plot above, change the value of that variable to something new and make a new prediction for that one case (i.e. set the value to something very small, or very large). How does the prediction change?

```
In [12]: 1 d['gavr'] = 40
          2 rf.predict_proba(d_reshaped)
```

```
C:\Users\Travis\AppData\Local\Temp\ipykernel_23268\2588936456.py:1: SettingWithCopyWarning:
```

```
A value is trying to be set on a copy of a slice from a DataFrame
```

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy) ([https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy))

```
d['gavr'] = 40
```

```
D:\Programming\Anaconda\lib\site-packages\sklearn\base.py:450: UserWarning: X does not have valid feature names, but RandomForestClassifier was fitted with feature names
```

```
warnings.warn(
```

```
Out[12]: array([[0.546, 0.454]])
```

Let's now try and predict the outcome for this case if that important variable was changed from it's minimum value to it's maximum value.

1. Create a grid of at least 100 points from the minimum value to the maximum value
2. Duplicate the case you used above the same number of times.
3. Add the grid of points to the data.
4. Predict the outcome using this new fake data and save the predicted probability in the dataset



```

In [19]: 1 # Step 1: Create a grid of at least 100 points from the minimum value to t
2 N = 100
3 g = np.linspace(np.min(X_train['gavr']), np.max(X_train['gavr']), num=N)
4
5 # Step 2: Duplicate the case you used above the same number of times.
6 X_train_dup = X_train.loc[[1] * 100, :]
7
8 # Step 3: Add the grid of points to the data.
9 X_train_dup.loc[:, 'gavr'] = g
10
11 # Step 4: Predict the outcome using this new fake data and save the predic
12 probs = rf.predict_proba(X_train_dup)[ :, 1]
13 X_train_dup['predicted_prob'] = probs
14 X_train_dup = X_train_dup.reset_index(drop=True)
15
16 #print(X_train_dup)
17 #print(probs)
18 #print(g)

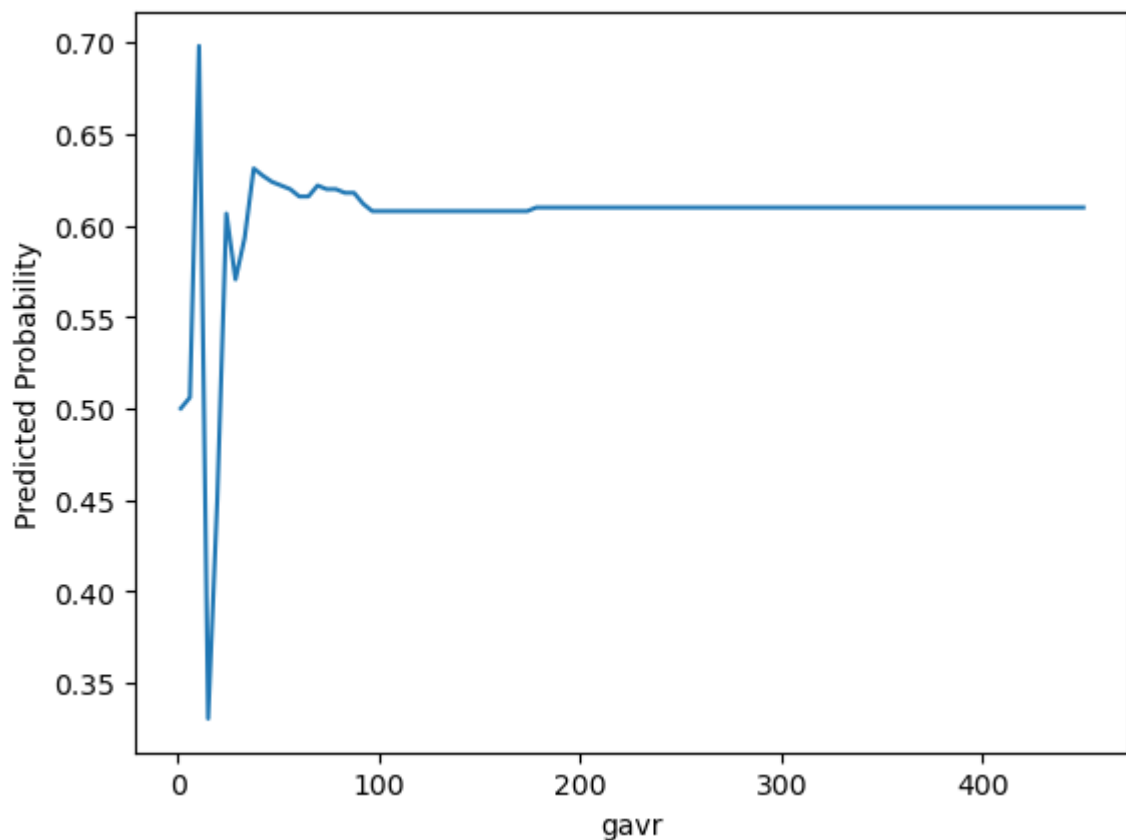
```

Now plot the results of that sequential grid against the predicted probability. How do you see the probability of responding the mailer change in response to the variable? Tip: You can use `coord_cartesian` to change the xlimits to focus on specific areas of detail if you want

```

In [14]: 1 #Plot the results
2 plt.plot(X_train_dup['gavr'], X_train_dup['predicted_prob'])
3 plt.xlabel('gavr')
4 plt.ylabel('Predicted Probability')
5 plt.show()

```



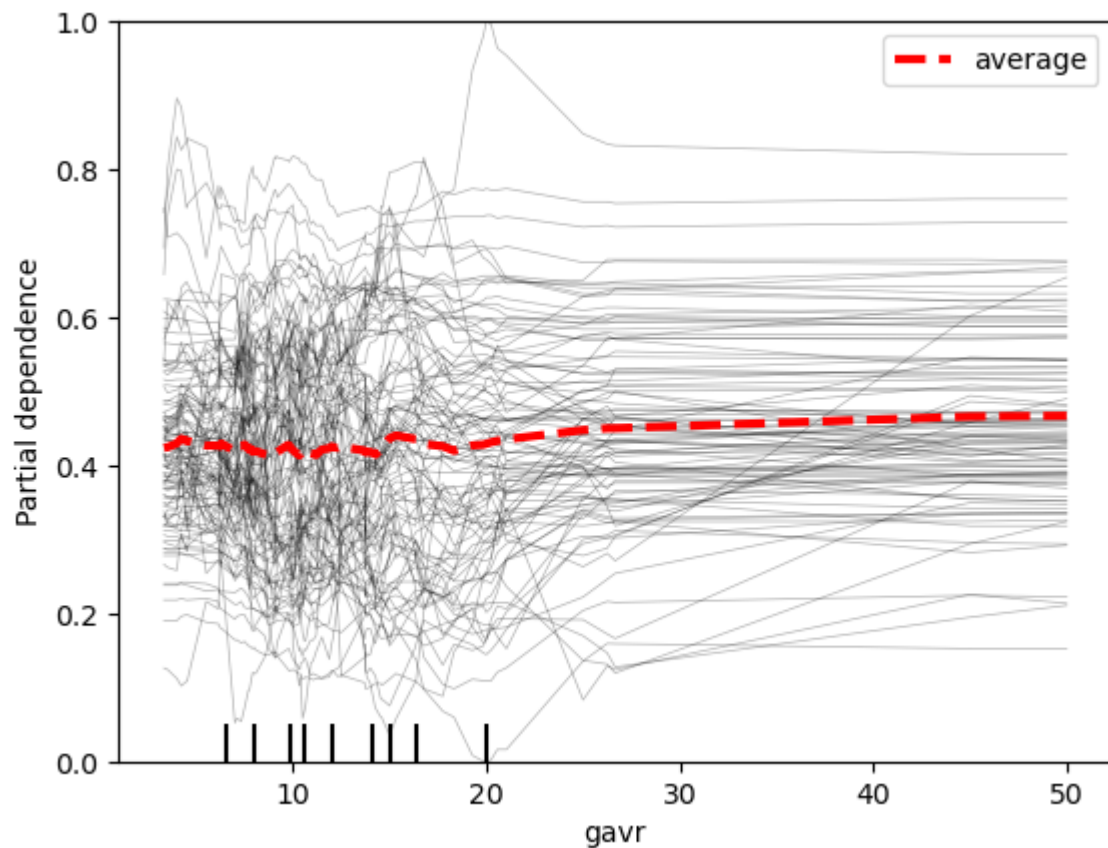
## Bonus: Use IML to create a PDP+ICE chart

Use the `iml` package to create an individual conditional expectations combined with partial dependence chart.

It's highly recommended that you only provide the prediction object a subset of the data. Use `coord_cartesian` to focus the y-limit range to focus the chart to a region where you can observe the effect (it's tiny!).

```
In [18]: 1 from sklearn.inspection import PartialDependenceDisplay
2 from sklearn.inspection import partial_dependence
3
4 X_test_100 = X_test.sample(n=100)
5 features = ["gavr"]
6 PartialDependenceDisplay.from_estimator(rf, X_test_100, features,
7                                         kind='both',
8                                         ice_lines_kw={"color": "black"},
9                                         pd_line_kw={"color": "red", "lw": 3,
```

```
Out[18]: <sklearn.inspection._plot.partial_dependence.PartialDependenceDisplay at 0x21400714310>
```



## 4. Summarize

***What are your thoughts on the impact of these different features on the likelihood for a person to respond to our donation requests? (100-200 words)***

Things you could write about:

1. What changes would you make to the modeling approach to improve the predictions?
2. What recommendations would you make to the organization?
3. What other kinds of data would be useful in improving these predictions?

### Summary

According to the code, the model is overfitting to the data. In order to reduce overfitting, it is recommended to regularize the random forest by finding a good mtry parameter value such that the test and train AUC are more similar, although it could impact training performance, the test performance can be improved.

In order to improve the performance of the model, we could consider the model selection. Using a different model other than the random forest model such as the K-nearest neighbours, logistic regression and support vector machines may occur a better result.

People are more likely to continue to donate if the average amount of gift is about 10.6. It is recommended to target people who have given an average amount of 8-12 for a higher chance of receiving a donation. Moreover, people who have an average donation of larger than 37 although has a slightly lower chance than 10.6, still has a moderate chance of giving out a donation. It is also recommended to target people who have an average donation of 37 since it is a relatively high amount and it would be more efficient to target the group as well.

Other kinds of data which would be useful to include are as follows:

Age: certain age groups may be more likely to donate to charities.

Status of the Receiver: Individuals who are married may be more likely to donate to charities.

Education level: Individuals with higher education level may be more likely to donate to charities.

Occupation: Individuals who work in healthcare places such as a hospital may be more likely to donate to charities.

Geographic location: Certain locations may be more charitable than others.