In [ ]:

In [1]:
```python
import pandas as pd

import numpy as np


from sklearn.model_selection import GridSearchCV

from sklearn.metrics import roc_curve,auc
from sklearn.model_selection import cross_val_score

from sklearn.metrics import accuracy_score
from sklearn.model_selection import cross_validate



from sklearn.metrics import accuracy_score,confusion_matrix,f1_score,precision_score,recall_score,roc_auc_score
from sklearn.model_selection import GridSearchCV


import xgboost as xgb
from sklearn.model_selection import RandomizedSearchCV
from xgboost import XGBRegressor

from sklearn.tree import DecisionTreeRegressor

from imblearn.over_sampling import SMOTE
from sklearn.model_selection import train_test_split

import warnings
warnings.filterwarnings('ignore')

from sklearn.preprocessing import StandardScaler

from sklearn.ensemble import RandomForestClassifier

from sklearn.metrics import roc_auc_score
from sklearn import metrics
from sklearn.metrics import accuracy_score,confusion_matrix,f1_score,precision_score,recall_score,roc_auc_score

import numpy as np;
import seaborn as sns; sns.set()
import matplotlib.pyplot as plt

from xgboost import XGBClassifier
from sklearn import linear_model
```

Using TensorFlow backend.

In [2]:
```python
df_tri=pd.read_csv('train.csv')
```

```
In [11]: df_tri.head(2)
```

Out[11]:

| | ID | Attr1 | Attr2 | Attr3 | Attr4 | Attr5 | Attr6 | Attr7 | Attr8 | Attr9 | ... | Attr |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0.135370 | 0.45185 | 0.31162 | 2.0469 | 10.2340 | 0.16768 | 0.167630 | 1.2131 | 2.2554 | ... | 0.0787 |
| 1 | 2 | 0.005861 | 0.39858 | 0.19768 | 1.9390 | 9.5771 | 0.00000 | 0.007237 | 1.5089 | 0.9788 | ... | 0.2697 |

2 rows × 66 columns

```
In [3]: df_tei=pd.read_csv('test.csv')
```

```
In [13]: df_tei.head(2)
```

Out[13]:

| | ID | Attr1 | Attr2 | Attr3 | Attr4 | Attr5 | Attr6 | Attr7 | Attr8 | Attr9 | ... | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 36554 | 0.20055 | 0.37951 | 0.396410 | 2.0472 | 32.351 | 0.38825 | 0.249760 | 1.33050 | 1.1389 | ... | 34 |
| 1 | 36555 | 0.00902 | 0.63202 | 0.053735 | 1.1263 | -37.842 | 0.00000 | 0.014434 | 0.58223 | 1.3332 | ... | |

2 rows × 65 columns

```
In [14]: df_tri['target'].value_counts()
```

```
Out[14]: 0    29772
         1     1511
         Name: target, dtype: int64
```

**OBSERVATION: As we can see here this is highly imbalance dataset:**

```
In [16]: yees=df_tri['target']
         df_tr_after_drop=df_tri.drop(['target'],axis=1)
```

# Splitting DATA into Train Test and C.V with stratification in 49:30:21 ratio:

```
In [17]: # split the data set into train and test
         X_1, X_test, y_1, y_test = train_test_split(df_tr_after_drop, yees, test_size=
         0.3, random_state=42,stratify=yees)

         # split the train data set into cross validation train and cross validation te
         st
         X_tr, X_cv, y_tr, y_cv = train_test_split(X_1, y_1, test_size=0.3)
```

## Replacing NaN with mean value of feature:

```
In [18]:  X_tr.fillna(X_tr.mean(), inplace=True)
          X_test.fillna(X_test.mean(), inplace=True)
          X_cv.fillna(X_cv.mean(), inplace=True)
```

## Standardisation of data:

```
In [19]:  sc = StandardScaler(with_mean=True)
          Xbow_tr_std = sc.fit_transform(X_tr)
          Xbow_test_std = sc.transform(X_test)
          Xbow_cv_std = sc.transform(X_cv)

          df_tei1 = sc.fit_transform(df_tei)
          df_tri1 = sc.transform(df_tr_after_drop)
```

```
C:\Users\all\AppData\Local\conda\conda\envs\tf\lib\site-packages\sklearn\prep
rocessing\data.py:645: DataConversionWarning: Data with input dtype int64, fl
oat64 were all converted to float64 by StandardScaler.
  return self.partial_fit(X, y)
C:\Users\all\AppData\Local\conda\conda\envs\tf\lib\site-packages\sklearn\bas
e.py:464: DataConversionWarning: Data with input dtype int64, float64 were al
l converted to float64 by StandardScaler.
  return self.fit(X, **fit_params).transform(X)
C:\Users\all\AppData\Local\conda\conda\envs\tf\lib\site-packages\ipykernel_la
uncher.py:3: DataConversionWarning: Data with input dtype int64, float64 were
all converted to float64 by StandardScaler.
  This is separate from the ipykernel package so we can avoid doing imports u
ntil
C:\Users\all\AppData\Local\conda\conda\envs\tf\lib\site-packages\ipykernel_la
uncher.py:4: DataConversionWarning: Data with input dtype int64, float64 were
all converted to float64 by StandardScaler.
  after removing the cwd from sys.path.
C:\Users\all\AppData\Local\conda\conda\envs\tf\lib\site-packages\sklearn\prep
rocessing\data.py:645: DataConversionWarning: Data with input dtype int64, fl
oat64 were all converted to float64 by StandardScaler.
  return self.partial_fit(X, y)
C:\Users\all\AppData\Local\conda\conda\envs\tf\lib\site-packages\sklearn\bas
e.py:464: DataConversionWarning: Data with input dtype int64, float64 were al
l converted to float64 by StandardScaler.
  return self.fit(X, **fit_params).transform(X)
C:\Users\all\AppData\Local\conda\conda\envs\tf\lib\site-packages\ipykernel_la
uncher.py:7: DataConversionWarning: Data with input dtype int64, float64 were
all converted to float64 by StandardScaler.
  import sys
```

# HYPER_PARAMETER TUNING WITH DEPTH AND NO.OF ESTIMATERS OF DECISION TREE:

```python
In [21]: #code for hyperparameter tuning
         import numpy as np
         hyper1 = [5, 10, 50, 100, 200,400, 600, 800]
         hyper2 = [2, 4, 6, 8, 10, 12, 14, 16, 18]

         auc1=np.empty((8,9))
         auc2=np.empty((8,9))
         l=0
         for j in hyper1:
             m=0
             for k in hyper2:

                 model = RandomForestClassifier(n_estimators=j,max_depth=k)
                 model.fit(Xbow_tr_std, y_tr)

                 probs = model.predict_proba(Xbow_tr_std)
                 preds = probs[:,1]
                 roc_auc1=metrics.roc_auc_score(y_tr, preds)
                 auc1[l][m]=(roc_auc1)


                 probs = model.predict_proba(Xbow_cv_std)
                 preds = probs[:,1]
                 roc_auc2=metrics.roc_auc_score(y_cv, preds)
                 auc2[l][m]=(roc_auc2)
                 m=m+1

             l=l+1
```
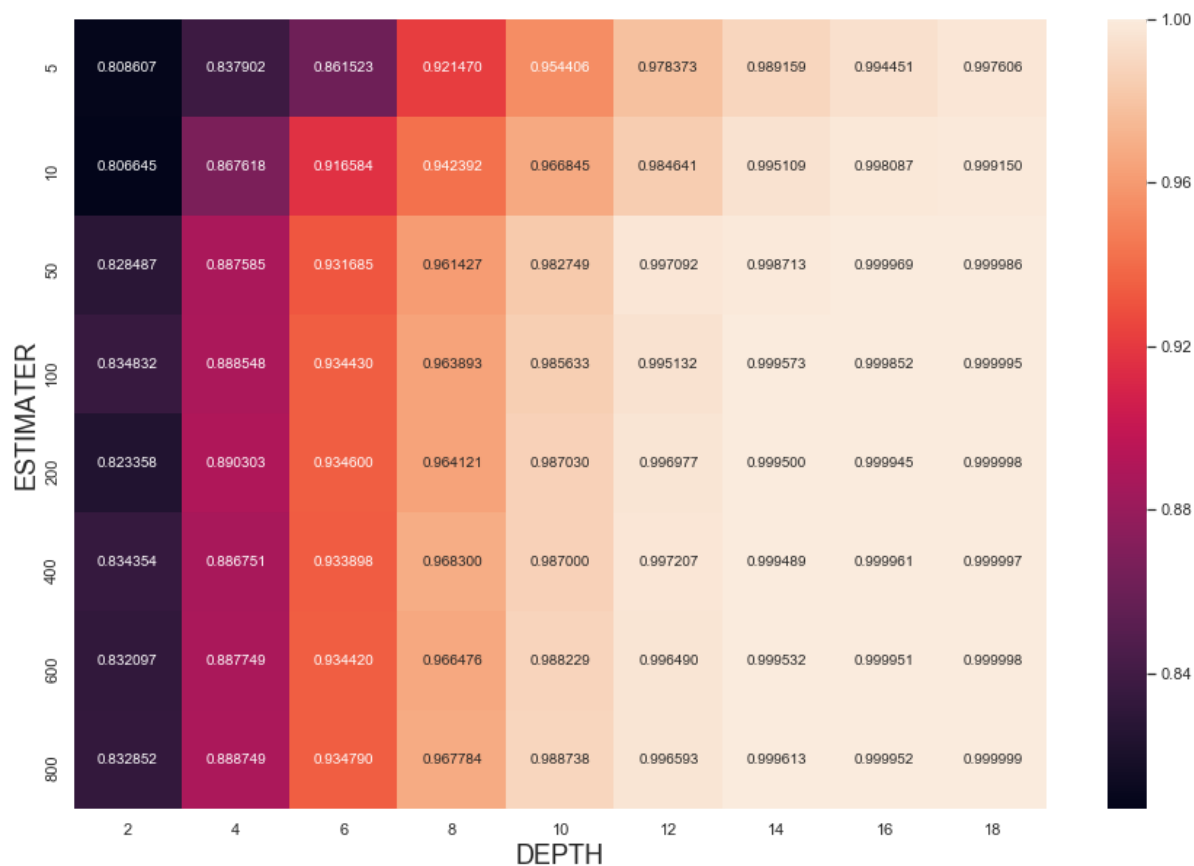
# PLOTTING SEABORN HEATMAP:

In [22]:
```python
# Code for drawing seaborn heatmaps

df_heatmap =pd.DataFrame(auc1, hyper1, hyper2 )
fig = plt.figure(figsize=(15,10))
ax = sns.heatmap(df_heatmap, annot=True, fmt="f")

plt.ylabel('ESTIMATER',size=18)
plt.xlabel('DEPTH',size=18)
plt.title("HEATMAP Matrix for train\n",size=24)
plt.show()
```

## HEATMAP Matrix for train

| ESTIMATER \ DEPTH | 2 | 4 | 6 | 8 | 10 | 12 | 14 | 16 | 18 |
|---|---|---|---|---|---|---|---|---|---|
| 5 | 0.808607 | 0.837902 | 0.861523 | 0.921470 | 0.954406 | 0.978373 | 0.989159 | 0.994451 | 0.997606 |
| 10 | 0.806645 | 0.867618 | 0.916584 | 0.942392 | 0.966845 | 0.984641 | 0.995109 | 0.998087 | 0.999150 |
| 50 | 0.828487 | 0.887585 | 0.931685 | 0.961427 | 0.982749 | 0.997092 | 0.998713 | 0.999969 | 0.999986 |
| 100 | 0.834832 | 0.888548 | 0.934430 | 0.963893 | 0.985633 | 0.995132 | 0.999573 | 0.999852 | 0.999995 |
| 200 | 0.823358 | 0.890303 | 0.934600 | 0.964121 | 0.987030 | 0.996977 | 0.999500 | 0.999945 | 0.999998 |
| 400 | 0.834354 | 0.886751 | 0.933898 | 0.968300 | 0.987000 | 0.997207 | 0.999489 | 0.999961 | 0.999997 |
| 600 | 0.832097 | 0.887749 | 0.934420 | 0.966476 | 0.988229 | 0.996490 | 0.999532 | 0.999951 | 0.999998 |
| 800 | 0.832852 | 0.888749 | 0.934790 | 0.967784 | 0.988738 | 0.996593 | 0.999613 | 0.999952 | 0.999999 |

In [23]:
```python
# Code for drawing seaborn heatmaps

df_heatmap =pd.DataFrame(auc2, hyper1, hyper2 )
fig = plt.figure(figsize=(15,10))
ax = sns.heatmap(df_heatmap, annot=True, fmt="f")

plt.ylabel('ESTIMATER',size=18)
plt.xlabel('DEPTH',size=18)
plt.title("HEATMAP Matrix for CV\n",size=24)
plt.show()
```
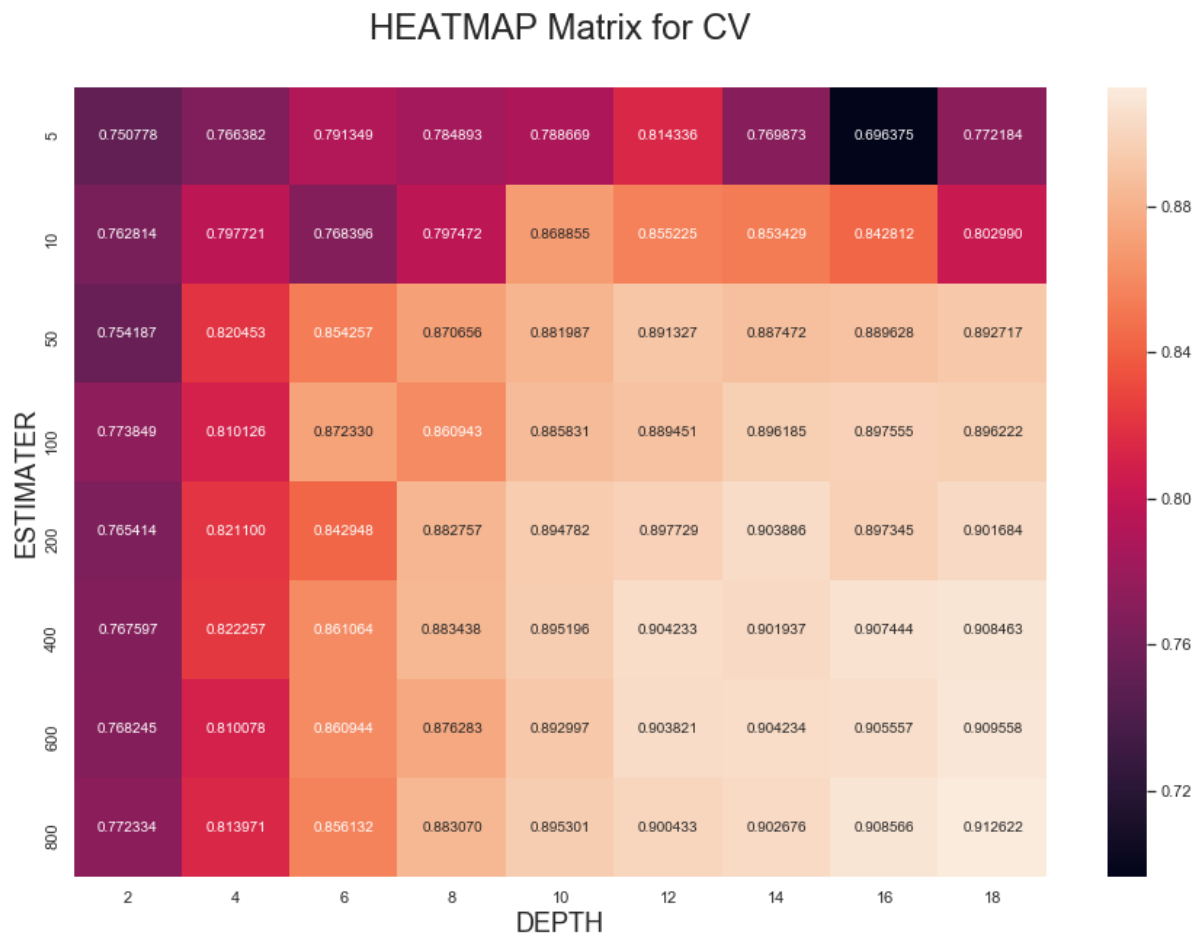
## HEATMAP Matrix for CV

| ESTIMATER \ DEPTH | 2 | 4 | 6 | 8 | 10 | 12 | 14 | 16 | 18 |
|---|---|---|---|---|---|---|---|---|---|
| 5 | 0.750778 | 0.766382 | 0.791349 | 0.784893 | 0.788669 | 0.814336 | 0.769873 | 0.696375 | 0.772184 |
| 10 | 0.762814 | 0.797721 | 0.768396 | 0.797472 | 0.868855 | 0.855225 | 0.853429 | 0.842812 | 0.802990 |
| 50 | 0.754187 | 0.820453 | 0.854257 | 0.870656 | 0.881987 | 0.891327 | 0.887472 | 0.889628 | 0.892717 |
| 100 | 0.773849 | 0.810126 | 0.872330 | 0.860943 | 0.885831 | 0.889451 | 0.896185 | 0.897555 | 0.896222 |
| 200 | 0.765414 | 0.821100 | 0.842948 | 0.882757 | 0.894782 | 0.897729 | 0.903886 | 0.897345 | 0.901684 |
| 400 | 0.767597 | 0.822257 | 0.861064 | 0.883438 | 0.895196 | 0.904233 | 0.901937 | 0.907444 | 0.908463 |
| 600 | 0.768245 | 0.810078 | 0.860944 | 0.876283 | 0.892997 | 0.903821 | 0.904234 | 0.905557 | 0.909558 |
| 800 | 0.772334 | 0.813971 | 0.856132 | 0.883070 | 0.895301 | 0.900433 | 0.902676 | 0.908566 | 0.912622 |

# FITTING AND TESTING MODEL ON OUR SPLITTED TEST DATA:

In [24]:
```python
rf = RandomForestClassifier(n_estimators=50,max_depth=8)

# fitting the model
rf.fit(Xbow_tr_std, y_tr)

# predict the response
pred = rf.predict(Xbow_test_std)

# evaluate accuracy
acc = accuracy_score(y_test, pred) * 100

precision_score1=precision_score(y_test, pred )

recall_score1=recall_score(y_test, pred )
f1 = f1_score(y_test, pred)


print('\nThe accuracy of the Random forest classifier for n_estimaters=%f and
 Depth = %f is %f%%' % (50,8, acc))



print('\nThe precision_score of the  Random forest classifier  for n_estimater
s=%d and Depth = %d is %f' % (50,8,precision_score1))

print('\nThe recall_score of the  Random forest classifier  for n_estimaters=%
d and Depth = %d is %f' % (50,8,recall_score1))

print('\nThe f1_score of the  Random forest classifier  for n_estimaters=%d an
d Depth = %d is %f' % (50,8,f1))
```

The accuracy of the Random forest classifier for n_estimaters=50.000000 and D
epth = 8.000000 is 95.279702%

The precision_score of the  Random forest classifier  for n_estimaters=50 and
Depth = 8 is 1.000000

The recall_score of the  Random forest classifier  for n_estimaters=50 and De
pth = 8 is 0.022075

The f1_score of the  Random forest classifier  for n_estimaters=50 and Depth
= 8 is 0.043197

# PLOTTING CONFUSION MATRIX:

In [25]:
```python
# Code for drawing seaborn heatmaps
class_names = ['negative','positive']
df_heatmap = pd.DataFrame(confusion_matrix(y_test, pred), index=class_names, c
olumns=class_names )
fig = plt.figure(figsize=(10,7))
heatmap = sns.heatmap(df_heatmap, annot=True, fmt="d")

# Setting tick labels for heatmap
heatmap.yaxis.set_ticklabels(heatmap.yaxis.get_ticklabels(), rotation=0, ha='r
ight', fontsize=14)
heatmap.xaxis.set_ticklabels(heatmap.xaxis.get_ticklabels(), rotation=0, ha='r
ight', fontsize=14)
plt.ylabel('Predicted label',size=18)
plt.xlabel('True label',size=18)
plt.title("Confusion Matrix\n",size=24)
plt.show()
```
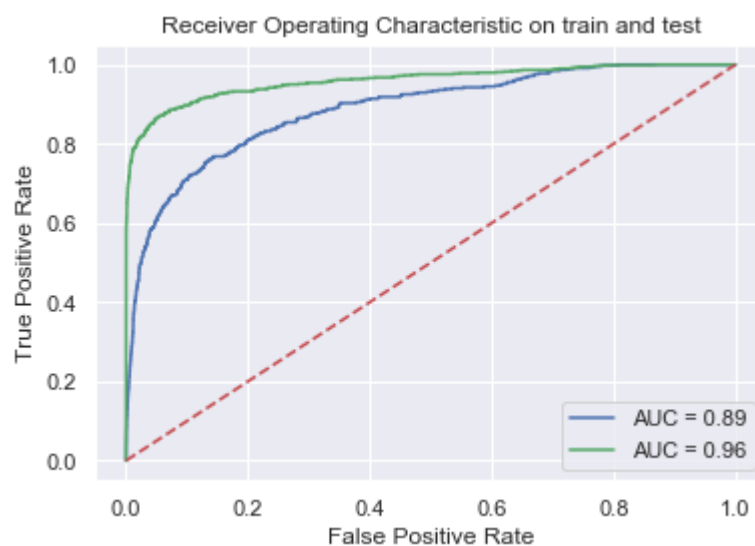
## Confusion Matrix



# PLOTTING AUC_ROC CURVE FOR TRAIN AND TEST DATA:

```
In [26]: rf.fit(Xbow_tr_std, y_tr)
         probs2 = rf.predict_proba(Xbow_tr_std)
         preds2 = probs2[:,1]
         fpr2, tpr2, threshold2 = metrics.roc_curve(y_tr, preds2)
         roc_auc2 = metrics.auc(fpr2, tpr2)


         probs1 = rf.predict_proba(Xbow_test_std)
         preds1 = probs1[:,1]
         fpr1, tpr1, threshold1 = metrics.roc_curve(y_test, preds1)
         roc_auc1 = metrics.auc(fpr1, tpr1)
```

```
In [27]: plt.title('Receiver Operating Characteristic on train and test')
         plt.plot(fpr1, tpr1, 'b', label = 'AUC = %0.2f' % roc_auc1)
         plt.plot(fpr2, tpr2, 'g', label = 'AUC = %0.2f' % roc_auc2)
         #plt.plot(neighbors, auc1,'g')
         #plt.plot(neighbors, auc2,'r')
         plt.legend(loc = 'lower right')
         plt.plot([0, 1], [0, 1],'r--')
         plt.ylabel('True Positive Rate')
         plt.xlabel('False Positive Rate')
         plt.show()
```



# HYPERPARAMETER TUNING WITH DEPTH AND NO.OF ESTIMATERS for XGBClassifier

```
In [28]: #code for hyperparameter tuning
         import numpy as np
         hyper1 =  [5, 10, 50, 100, 200,400, 600, 800]
         hyper2 =  [2, 4, 6, 8, 10, 12, 14, 16, 18]

         auc1=np.empty((8,9))
         auc2=np.empty((8,9))
         l=0
         for j in hyper1:
             m=0
             for k in hyper2:

                 model = XGBClassifier(n_estimators=j,max_depth=k)
                 model.fit(Xbow_tr_std, y_tr)

                 probs = model.predict_proba(Xbow_tr_std)
                 preds = probs[:,1]
                 roc_auc1=metrics.roc_auc_score(y_tr, preds)
                 auc1[l][m]=(roc_auc1)


                 probs = model.predict_proba(Xbow_cv_std)
                 preds = probs[:,1]
                 roc_auc2=metrics.roc_auc_score(y_cv, preds)
                 auc2[l][m]=(roc_auc2)
                 m=m+1


             l=l+1
```

# PLOTTING SEABORN HEATMAP:

In [29]:
```python
# Code for drawing seaborn heatmaps

df_heatmap =pd.DataFrame(auc1, hyper1, hyper2 )
fig = plt.figure(figsize=(15,10))
ax = sns.heatmap(df_heatmap, annot=True, fmt="f")

plt.ylabel('ESTIMATER',size=18)
plt.xlabel('DEPTH',size=18)
plt.title("HEATMAP Matrix for train\n",size=24)
plt.show()
```
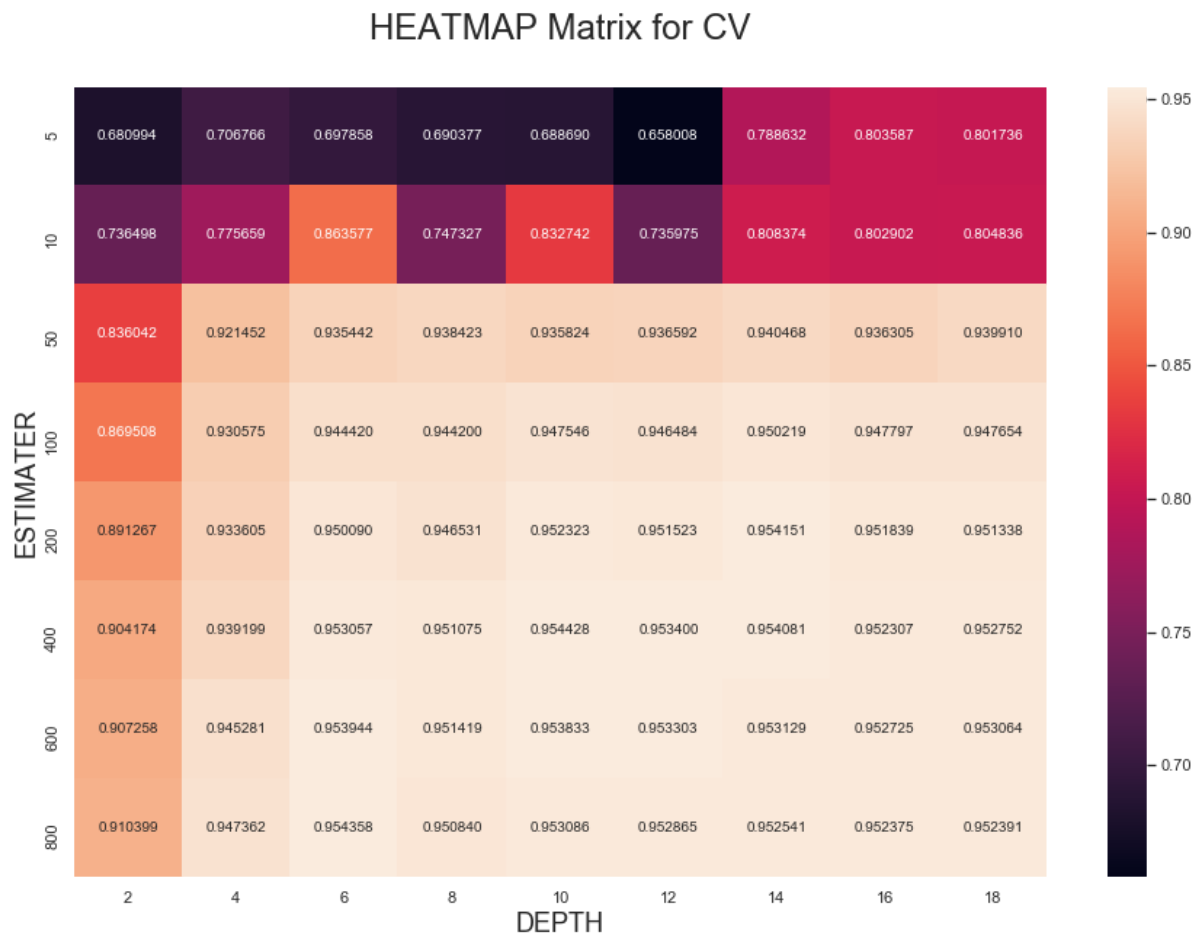
## HEATMAP Matrix for train

| | 2 | 4 | 6 | 8 | 10 | 12 | 14 | 16 | 18 |
|---|---|---|---|---|---|---|---|---|---|
| 5 | 0.803319 | 0.874229 | 0.920542 | 0.941859 | 0.952532 | 0.957489 | 0.966158 | 0.972904 | 0.972750 |
| 10 | 0.823787 | 0.900175 | 0.932545 | 0.953728 | 0.967987 | 0.973535 | 0.984319 | 0.985103 | 0.985342 |
| 50 | 0.888450 | 0.967648 | 0.993199 | 0.998850 | 0.999958 | 0.999996 | 0.999999 | 0.999999 | 0.999998 |
| 100 | 0.919561 | 0.988445 | 0.999356 | 0.999995 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 200 | 0.952955 | 0.998462 | 0.999999 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 400 | 0.978087 | 0.999996 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 600 | 0.988235 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 800 | 0.993416 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |

ESTIMATER / DEPTH

In [30]:
```python
# Code for drawing seaborn heatmaps

df_heatmap =pd.DataFrame(auc2, hyper1, hyper2 )
fig = plt.figure(figsize=(15,10))
ax = sns.heatmap(df_heatmap, annot=True, fmt="f")

plt.ylabel('ESTIMATER',size=18)
plt.xlabel('DEPTH',size=18)
plt.title("HEATMAP Matrix for CV\n",size=24)
plt.show()
```

## HEATMAP Matrix for CV

| ESTIMATER \ DEPTH | 2 | 4 | 6 | 8 | 10 | 12 | 14 | 16 | 18 |
|---|---|---|---|---|---|---|---|---|---|
| 5 | 0.680994 | 0.706766 | 0.697858 | 0.690377 | 0.688690 | 0.658008 | 0.788632 | 0.803587 | 0.801736 |
| 10 | 0.736498 | 0.775659 | 0.863577 | 0.747327 | 0.832742 | 0.735975 | 0.808374 | 0.802902 | 0.804836 |
| 50 | 0.836042 | 0.921452 | 0.935442 | 0.938423 | 0.935824 | 0.936592 | 0.940468 | 0.936305 | 0.939910 |
| 100 | 0.869508 | 0.930575 | 0.944420 | 0.944200 | 0.947546 | 0.946484 | 0.950219 | 0.947797 | 0.947654 |
| 200 | 0.891267 | 0.933605 | 0.950090 | 0.946531 | 0.952323 | 0.951523 | 0.954151 | 0.951839 | 0.951338 |
| 400 | 0.904174 | 0.939199 | 0.953057 | 0.951075 | 0.954428 | 0.953400 | 0.954081 | 0.952307 | 0.952752 |
| 600 | 0.907258 | 0.945281 | 0.953944 | 0.951419 | 0.953833 | 0.953303 | 0.953129 | 0.952725 | 0.953064 |
| 800 | 0.910399 | 0.947362 | 0.954358 | 0.950840 | 0.953086 | 0.952865 | 0.952541 | 0.952375 | 0.952391 |

# FITTING AND TESTING MODEL ON OUR SPLITTED TEST DATA:

In [34]:
```python
rf = RandomForestClassifier(n_estimators=50,max_depth=8)

# fitting the model
rf.fit(Xbow_tr_std, y_tr)

# predict the response
pred = rf.predict(Xbow_test_std)

# evaluate accuracy
acc = accuracy_score(y_test, pred) * 100

precision_score1=precision_score(y_test, pred )

recall_score1=recall_score(y_test, pred )

f1 = f1_score(y_test, pred)

print('\nThe accuracy of the Random forest classifier for n_estimaters=%f and
 Depth = %f is %f%%' % (50,8, acc))


print('\nThe precision_score of the  Random forest classifier  for n_estimater
s=%d and Depth = %d is %f' % (50,8,precision_score1))

print('\nThe recall_score of the  Random forest classifier  for n_estimaters=%
d and Depth = %d is %f' % (50,8,recall_score1))

print('\nThe f1_score of the  Random forest classifier  for n_estimaters=%d an
d Depth = %d is %f' % (50,8,f1))
```

The accuracy of the Random forest classifier for n_estimaters=50.000000 and D
epth = 8.000000 is 95.364944%

The precision_score of the  Random forest classifier  for n_estimaters=50 and
Depth = 8 is 0.950000

The recall_score of the  Random forest classifier  for n_estimaters=50 and De
pth = 8 is 0.041943

The f1_score of the  Random forest classifier  for n_estimaters=50 and Depth
= 8 is 0.080338

# PLOTTING CONFUSION MATRIX:

In [35]:
```python
# Code for drawing seaborn heatmaps
class_names = ['negative','positive']
df_heatmap = pd.DataFrame(confusion_matrix(y_test, pred), index=class_names, c
olumns=class_names )
fig = plt.figure(figsize=(10,7))
heatmap = sns.heatmap(df_heatmap, annot=True, fmt="d")

# Setting tick labels for heatmap
heatmap.yaxis.set_ticklabels(heatmap.yaxis.get_ticklabels(), rotation=0, ha='r
ight', fontsize=14)
heatmap.xaxis.set_ticklabels(heatmap.xaxis.get_ticklabels(), rotation=0, ha='r
ight', fontsize=14)
plt.ylabel('Predicted label',size=18)
plt.xlabel('True label',size=18)
plt.title("Confusion Matrix\n",size=24)
plt.show()
```



# PLOTTING AUC_ROC CURVE FOR TRAIN AND TEST DATA

In [36]:
```python
rf.fit(Xbow_tr_std, y_tr)
probs2 = rf.predict_proba(Xbow_tr_std)
preds2 = probs2[:,1]
fpr2, tpr2, threshold2 = metrics.roc_curve(y_tr, preds2)
roc_auc2 = metrics.auc(fpr2, tpr2)


probs1 = rf.predict_proba(Xbow_test_std)
preds1 = probs1[:,1]
fpr1, tpr1, threshold1 = metrics.roc_curve(y_test, preds1)
roc_auc1 = metrics.auc(fpr1, tpr1)
```

In [37]:
```python
plt.title('Receiver Operating Characteristic on train and test')
plt.plot(fpr1, tpr1, 'b', label = 'AUC = %0.2f' % roc_auc1)
plt.plot(fpr2, tpr2, 'g', label = 'AUC = %0.2f' % roc_auc2)
#plt.plot(neighbors, auc1,'g')
#plt.plot(neighbors, auc2,'r')
plt.legend(loc = 'lower right')
plt.plot([0, 1], [0, 1],'r--')
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()
```



# SMOTE

```
In [4]: df_tr=pd.read_csv('train.csv')
        yees=df_tr['target']
        df_tr_after_drop=df_tr.drop(['target'],axis=1)
        #spilliting data


        # split the data set into train and test
        X_1, X_test, y_1, y_test = train_test_split(df_tr_after_drop, yees, test_size=
        0.3, random_state=42,stratify=yees)

        # split the train data set into cross validation train and cross validation te
        st
        X_tr, X_cv, y_tr, y_cv = train_test_split(X_1, y_1, test_size=0.3)

        X_tr.fillna(X_tr.mean(), inplace=True)
        X_test.fillna(X_test.mean(), inplace=True)
        X_cv.fillna(X_cv.mean(), inplace=True)


        sm = SMOTE(random_state=27, ratio=1.0)

        X_tr, y_tr = sm.fit_sample(X_tr, y_tr)

        sc = StandardScaler(with_mean=True)
        Xbow_tr_std = sc.fit_transform(X_tr)
        Xbow_test_std = sc.transform(X_test)
        Xbow_cv_std = sc.transform(X_cv)
```

## HYPER_PARAMETER TUNING WITH DEPTH AND NO.OF ESTIMATERS OF DECISION TREE:

```
In [5]:  #code for hyperparameter tuning
         import numpy as np
         hyper1 =  [5, 10, 50, 100, 200,400, 600, 800]
         hyper2 =  [2, 4, 6, 8, 10, 12, 14, 16, 18]

         auc1=np.empty((8,9))
         auc2=np.empty((8,9))
         l=0
         for j in hyper1:
             m=0
             for k in hyper2:

                 model = RandomForestClassifier(n_estimators=j,max_depth=k)
                 model.fit(Xbow_tr_std, y_tr)

                 probs = model.predict_proba(Xbow_tr_std)
                 preds = probs[:,1]
                 roc_auc1=metrics.roc_auc_score(y_tr, preds)
                 auc1[l][m]=(roc_auc1)


                 probs = model.predict_proba(Xbow_cv_std)
                 preds = probs[:,1]
                 roc_auc2=metrics.roc_auc_score(y_cv, preds)
                 auc2[l][m]=(roc_auc2)
                 m=m+1

             l=l+1
```
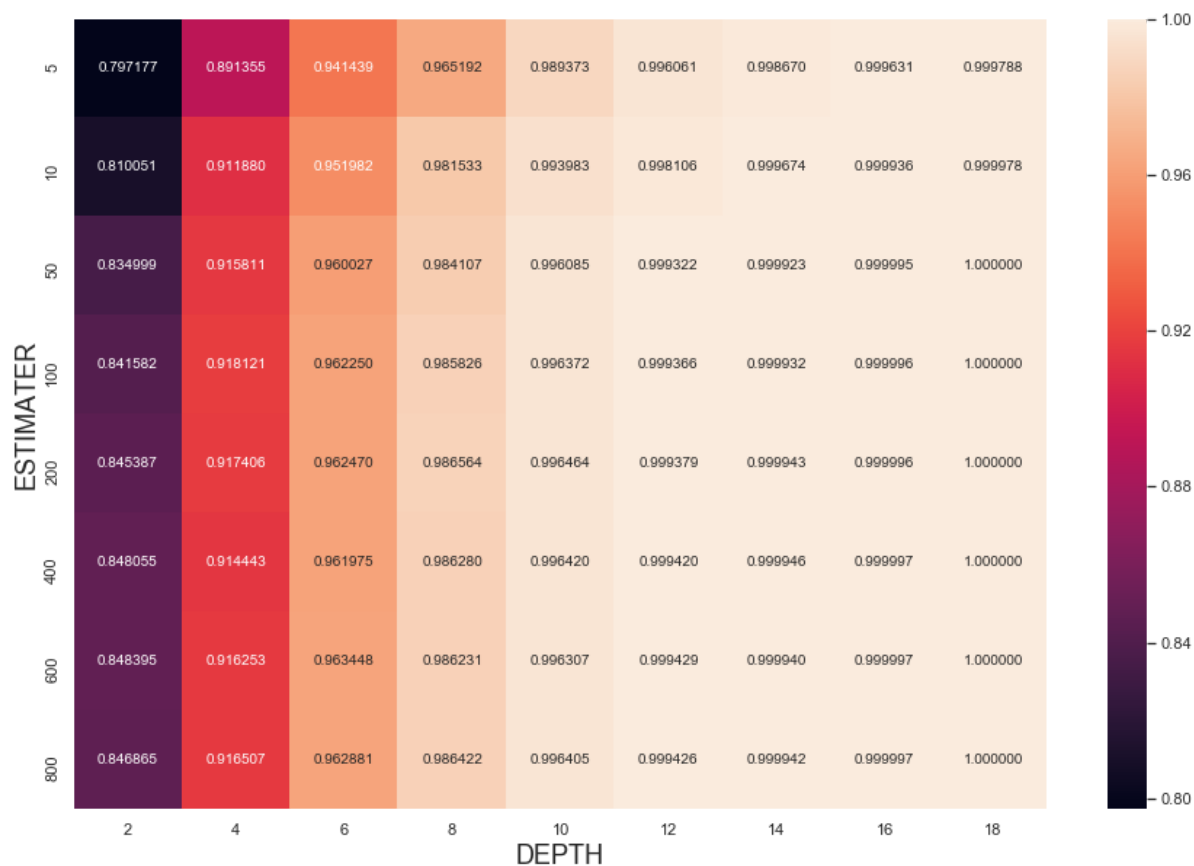
# PLOTTING SEABORN HEATMAP:

In [6]:
```python
# Code for drawing seaborn heatmaps

df_heatmap =pd.DataFrame(auc1, hyper1, hyper2 )
fig = plt.figure(figsize=(15,10))
ax = sns.heatmap(df_heatmap, annot=True, fmt="f")

plt.ylabel('ESTIMATER',size=18)
plt.xlabel('DEPTH',size=18)
plt.title("HEATMAP Matrix for train\n",size=24)
plt.show()
```
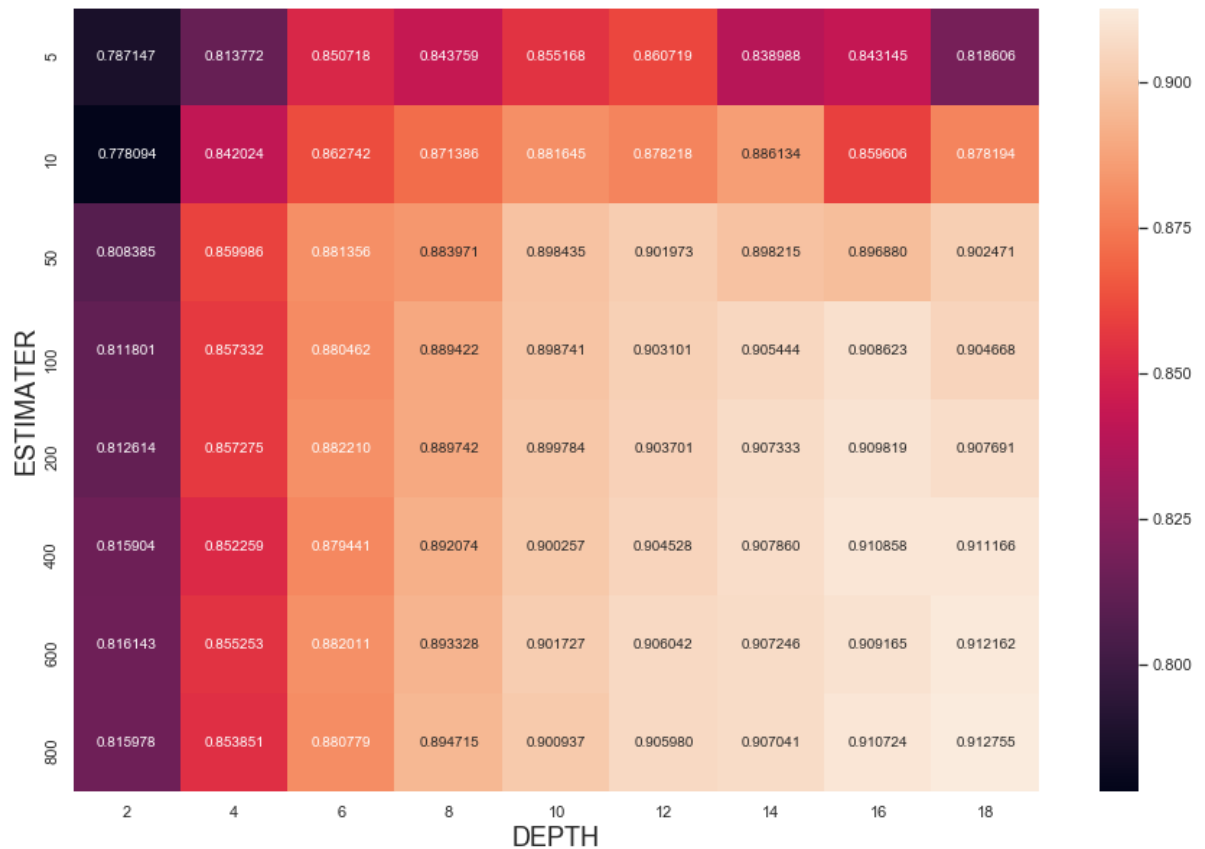
## HEATMAP Matrix for train

| ESTIMATER \ DEPTH | 2 | 4 | 6 | 8 | 10 | 12 | 14 | 16 | 18 |
|---|---|---|---|---|---|---|---|---|---|
| 5 | 0.797177 | 0.891355 | 0.941439 | 0.965192 | 0.989373 | 0.996061 | 0.998670 | 0.999631 | 0.999788 |
| 10 | 0.810051 | 0.911880 | 0.951982 | 0.981533 | 0.993983 | 0.998106 | 0.999674 | 0.999936 | 0.999978 |
| 50 | 0.834999 | 0.915811 | 0.960027 | 0.984107 | 0.996085 | 0.999322 | 0.999923 | 0.999995 | 1.000000 |
| 100 | 0.841582 | 0.918121 | 0.962250 | 0.985826 | 0.996372 | 0.999366 | 0.999932 | 0.999996 | 1.000000 |
| 200 | 0.845387 | 0.917406 | 0.962470 | 0.986564 | 0.996464 | 0.999379 | 0.999943 | 0.999996 | 1.000000 |
| 400 | 0.848055 | 0.914443 | 0.961975 | 0.986280 | 0.996420 | 0.999420 | 0.999946 | 0.999997 | 1.000000 |
| 600 | 0.848395 | 0.916253 | 0.963448 | 0.986231 | 0.996307 | 0.999429 | 0.999940 | 0.999997 | 1.000000 |
| 800 | 0.846865 | 0.916507 | 0.962881 | 0.986422 | 0.996405 | 0.999426 | 0.999942 | 0.999997 | 1.000000 |

In [7]:
```python
# Code for drawing seaborn heatmaps

df_heatmap =pd.DataFrame(auc2, hyper1, hyper2 )
fig = plt.figure(figsize=(15,10))
ax = sns.heatmap(df_heatmap, annot=True, fmt="f")

plt.ylabel('ESTIMATER',size=18)
plt.xlabel('DEPTH',size=18)
plt.title("HEATMAP Matrix for CV\n",size=24)
plt.show()
```

## HEATMAP Matrix for CV

| ESTIMATER \ DEPTH | 2 | 4 | 6 | 8 | 10 | 12 | 14 | 16 | 18 |
|---|---|---|---|---|---|---|---|---|---|
| 5 | 0.787147 | 0.813772 | 0.850718 | 0.843759 | 0.855168 | 0.860719 | 0.838988 | 0.843145 | 0.818606 |
| 10 | 0.778094 | 0.842024 | 0.862742 | 0.871386 | 0.881645 | 0.878218 | 0.886134 | 0.859606 | 0.878194 |
| 50 | 0.808385 | 0.859986 | 0.881356 | 0.883971 | 0.898435 | 0.901973 | 0.898215 | 0.896880 | 0.902471 |
| 100 | 0.811801 | 0.857332 | 0.880462 | 0.889422 | 0.898741 | 0.903101 | 0.905444 | 0.908623 | 0.904668 |
| 200 | 0.812614 | 0.857275 | 0.882210 | 0.889742 | 0.899784 | 0.903701 | 0.907333 | 0.909819 | 0.907691 |
| 400 | 0.815904 | 0.852259 | 0.879441 | 0.892074 | 0.900257 | 0.904528 | 0.907860 | 0.910858 | 0.911166 |
| 600 | 0.816143 | 0.855253 | 0.882011 | 0.893328 | 0.901727 | 0.906042 | 0.907246 | 0.909165 | 0.912162 |
| 800 | 0.815978 | 0.853851 | 0.880779 | 0.894715 | 0.900937 | 0.905980 | 0.907041 | 0.910724 | 0.912755 |

# FITTING AND TESTING MODEL ON OUR SPLITTED TEST DATA:

In [8]:
```python
rf = RandomForestClassifier(n_estimators=50,max_depth=10)

# fitting the model
rf.fit(Xbow_tr_std, y_tr)

# predict the response
pred = rf.predict(Xbow_test_std)

# evaluate accuracy
acc = accuracy_score(y_test, pred) * 100

precision_score1=precision_score(y_test, pred )

recall_score1=recall_score(y_test, pred )

print('\nThe accuracy of the Random forest classifier for n_estimaters=%f and
 Depth = %f is %f%%' % (50,10, acc))


f1 = f1_score(y_test, pred)
print('\nThe precision_score of the  Random forest classifier  for n_estimater
s=%d and Depth = %d is %f' % (50,10,precision_score1))

print('\nThe recall_score of the  Random forest classifier  for n_estimaters=%
d and Depth = %d is %f' % (50,10,recall_score1))

print('\nThe f1_score of the  Random forest classifier  for n_estimaters=%d an
d Depth = %d is %f' % (50,10,f1))
```

```
The accuracy of the Random forest classifier for n_estimaters=50.000000 and D
epth = 10.000000 is 90.602025%

The precision_score of the  Random forest classifier  for n_estimaters=50 and
Depth = 10 is 0.290323

The recall_score of the  Random forest classifier  for n_estimaters=50 and De
pth = 10 is 0.655629

The f1_score of the  Random forest classifier  for n_estimaters=50 and Depth
= 10 is 0.402439
```

# PLOTTING CONFUSION MATRIX:

In [9]:
```python
# Code for drawing seaborn heatmaps
class_names = ['negative','positive']
df_heatmap = pd.DataFrame(confusion_matrix(y_test, pred), index=class_names, c
olumns=class_names )
fig = plt.figure(figsize=(10,7))
heatmap = sns.heatmap(df_heatmap, annot=True, fmt="d")

# Setting tick labels for heatmap
heatmap.yaxis.set_ticklabels(heatmap.yaxis.get_ticklabels(), rotation=0, ha='r
ight', fontsize=14)
heatmap.xaxis.set_ticklabels(heatmap.xaxis.get_ticklabels(), rotation=0, ha='r
ight', fontsize=14)
plt.ylabel('Predicted label',size=18)
plt.xlabel('True label',size=18)
plt.title("Confusion Matrix\n",size=24)
plt.show()
```



# PLOTTING AUC_ROC CURVE FOR TRAIN AND TEST DATA:

```
In [10]: rf.fit(Xbow_tr_std, y_tr)
         probs2 = rf.predict_proba(Xbow_tr_std)
         preds2 = probs2[:,1]
         fpr2, tpr2, threshold2 = metrics.roc_curve(y_tr, preds2)
         roc_auc2 = metrics.auc(fpr2, tpr2)


         probs1 = rf.predict_proba(Xbow_test_std)
         preds1 = probs1[:,1]
         fpr1, tpr1, threshold1 = metrics.roc_curve(y_test, preds1)
         roc_auc1 = metrics.auc(fpr1, tpr1)
```

```
In [11]: plt.title('Receiver Operating Characteristic on train and test')
         plt.plot(fpr1, tpr1, 'b', label = 'AUC = %0.2f' % roc_auc1)
         plt.plot(fpr2, tpr2, 'g', label = 'AUC = %0.2f' % roc_auc2)
         #plt.plot(neighbors, auc1,'g')
         #plt.plot(neighbors, auc2,'r')
         plt.legend(loc = 'lower right')
         plt.plot([0, 1], [0, 1],'r--')
         plt.ylabel('True Positive Rate')
         plt.xlabel('False Positive Rate')
         plt.show()
```



# HYPERPARAMETER TUNING WITH DEPTH AND NO.OF ESTIMATERS for XGBClassifier

In [12]:
```python
#code for hyperparameter tuning
import numpy as np
hyper1 = [5, 10, 50, 100, 200,400, 600, 800]
hyper2 = [2, 4, 6, 8, 10, 12, 14, 16, 18]

auc1=np.empty((8,9))
auc2=np.empty((8,9))

l=0
for j in hyper1:
    m=0
    for k in hyper2:

        model = XGBClassifier(n_estimators=j,max_depth=k)
        model.fit(Xbow_tr_std, y_tr)

        probs = model.predict_proba(Xbow_tr_std)
        preds = probs[:,1]
        roc_auc1=metrics.roc_auc_score(y_tr, preds)
        auc1[l][m]=(roc_auc1)


        probs = model.predict_proba(Xbow_cv_std)
        preds = probs[:,1]
        roc_auc2=metrics.roc_auc_score(y_cv, preds)
        auc2[l][m]=(roc_auc2)
        m=m+1

    l=l+1
```
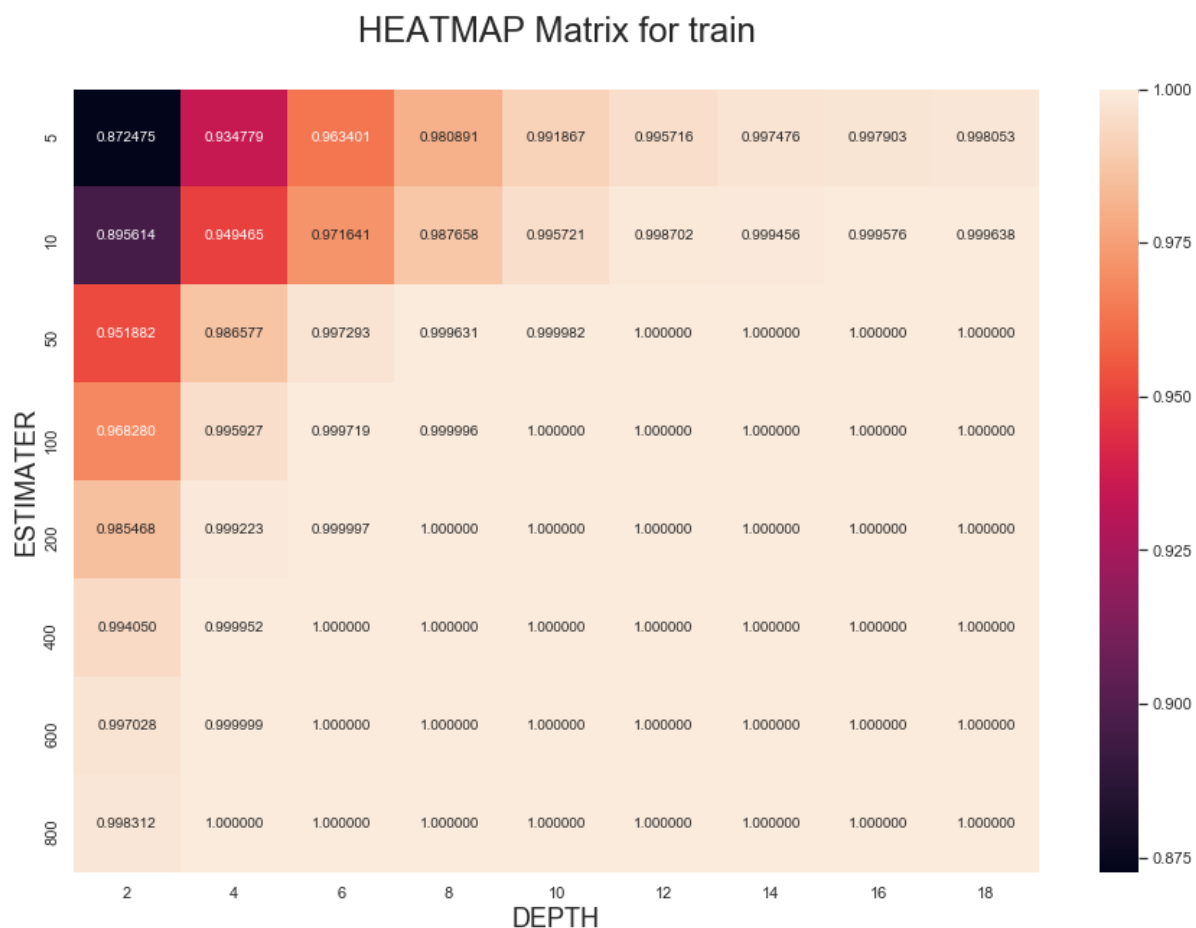
# PLOTTING SEABORN HEATMAP:

In [13]:
```python
# Code for drawing seaborn heatmaps

df_heatmap =pd.DataFrame(auc1, hyper1, hyper2 )
fig = plt.figure(figsize=(15,10))
ax = sns.heatmap(df_heatmap, annot=True, fmt="f")

plt.ylabel('ESTIMATER',size=18)
plt.xlabel('DEPTH',size=18)
plt.title("HEATMAP Matrix for train\n",size=24)
plt.show()
```
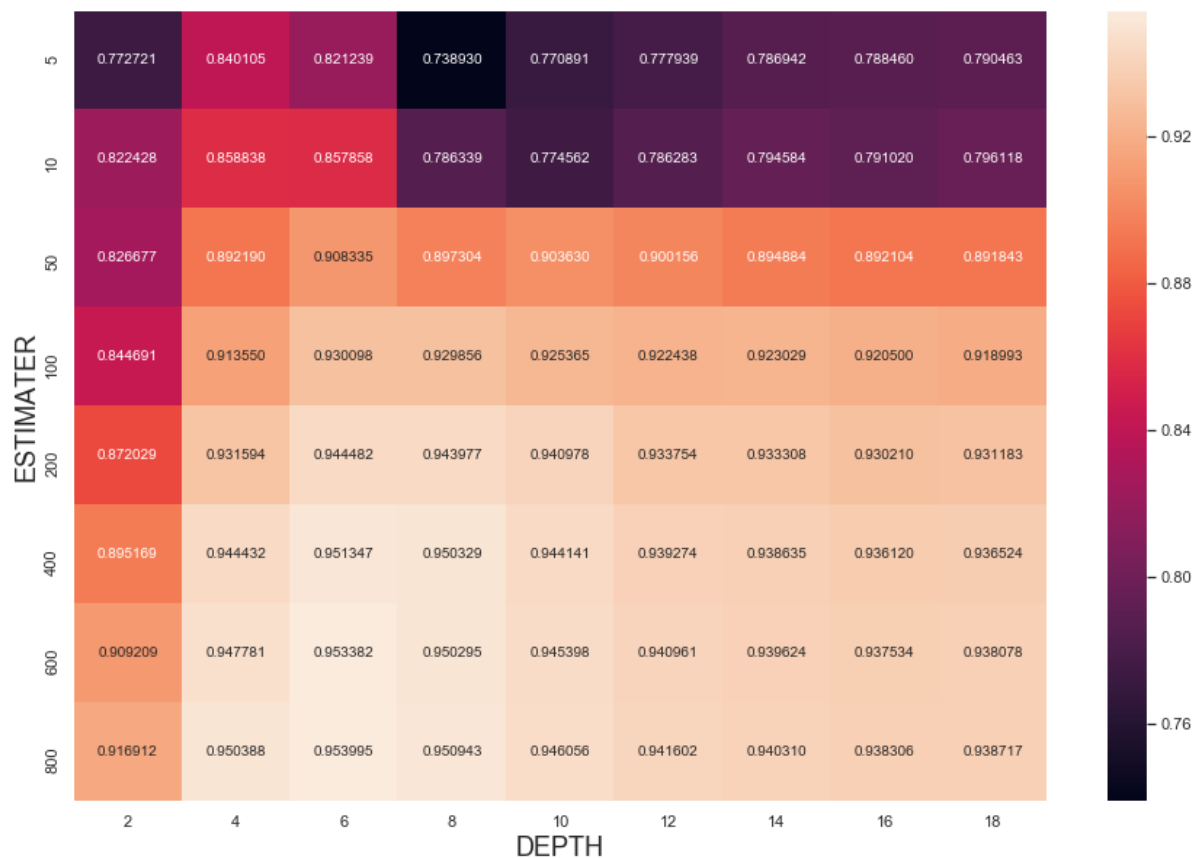
## HEATMAP Matrix for train

| ESTIMATER \ DEPTH | 2 | 4 | 6 | 8 | 10 | 12 | 14 | 16 | 18 |
|---|---|---|---|---|---|---|---|---|---|
| 5 | 0.872475 | 0.934779 | 0.963401 | 0.980891 | 0.991867 | 0.995716 | 0.997476 | 0.997903 | 0.998053 |
| 10 | 0.895614 | 0.949465 | 0.971641 | 0.987658 | 0.995721 | 0.998702 | 0.999456 | 0.999576 | 0.999638 |
| 50 | 0.951882 | 0.986577 | 0.997293 | 0.999631 | 0.999982 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 100 | 0.968280 | 0.995927 | 0.999719 | 0.999996 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 200 | 0.985468 | 0.999223 | 0.999997 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 400 | 0.994050 | 0.999952 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 600 | 0.997028 | 0.999999 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 800 | 0.998312 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |

In [14]:
```python
# Code for drawing seaborn heatmaps

df_heatmap =pd.DataFrame(auc2, hyper1, hyper2 )
fig = plt.figure(figsize=(15,10))
ax = sns.heatmap(df_heatmap, annot=True, fmt="f")

plt.ylabel('ESTIMATER',size=18)
plt.xlabel('DEPTH',size=18)
plt.title("HEATMAP Matrix for CV\n",size=24)
plt.show()
```



HEATMAP Matrix for CV

# FITTING AND TESTING MODEL ON OUR SPLITTED TEST DATA:

In [15]:
```python
rf = RandomForestClassifier(n_estimators=100,max_depth=8)

# fitting the model
rf.fit(Xbow_tr_std, y_tr)

# predict the response
pred = rf.predict(Xbow_test_std)

# evaluate accuracy
acc = accuracy_score(y_test, pred) * 100

precision_score1=precision_score(y_test, pred )

recall_score1=recall_score(y_test, pred )

f1 = f1_score(y_test, pred)

print('\nThe accuracy of the Random forest classifier for n_estimaters=%f and
 Depth = %f is %f%%' % (100,8, acc))



print('\nThe precision_score of the  Random forest classifier  for n_estimater
s=%d and Depth = %d is %f' % (100,8,precision_score1))

print('\nThe recall_score of the  Random forest classifier  for n_estimaters=%
d and Depth = %d is %f' % (100,8,recall_score1))

print('\nThe f1_score of the  Random forest classifier  for n_estimaters=%d an
d Depth = %d is %f' % (100,8,f1))
```

The accuracy of the Random forest classifier for n_estimaters=100.000000 and
Depth = 8.000000 is 88.609483%

The precision_score of the  Random forest classifier  for n_estimaters=100 an
d Depth = 8 is 0.248366

The recall_score of the  Random forest classifier  for n_estimaters=100 and D
epth = 8 is 0.671082

The f1_score of the  Random forest classifier  for n_estimaters=100 and Depth
= 8 is 0.362552

# PLOTTING CONFUSION MATRIX:

```
In [16]:  # Code for drawing seaborn heatmaps
          class_names = ['negative','positive']
          df_heatmap = pd.DataFrame(confusion_matrix(y_test, pred), index=class_names, c
          olumns=class_names )
          fig = plt.figure(figsize=(10,7))
          heatmap = sns.heatmap(df_heatmap, annot=True, fmt="d")

          # Setting tick labels for heatmap
          heatmap.yaxis.set_ticklabels(heatmap.yaxis.get_ticklabels(), rotation=0, ha='r
          ight', fontsize=14)
          heatmap.xaxis.set_ticklabels(heatmap.xaxis.get_ticklabels(), rotation=0, ha='r
          ight', fontsize=14)
          plt.ylabel('Predicted label',size=18)
          plt.xlabel('True label',size=18)
          plt.title("Confusion Matrix\n",size=24)
          plt.show()
```
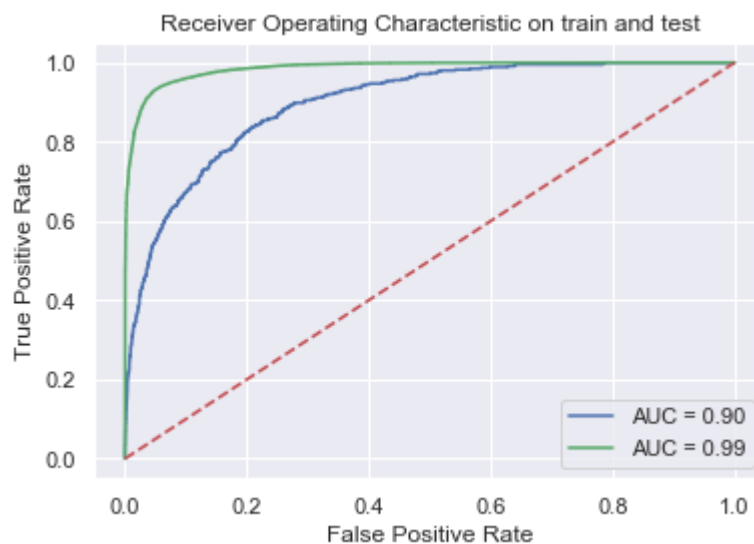
## Confusion Matrix

| | negative | positive |
|---|---|---|
| negative | 8012 | 920 |
| positive | 149 | 304 |

Predicted label (y-axis), True label (x-axis)

# PLOTTING AUC_ROC CURVE FOR TRAIN AND TEST DATA:

In [17]:
```
rf.fit(Xbow_tr_std, y_tr)
probs2 = rf.predict_proba(Xbow_tr_std)
preds2 = probs2[:,1]
fpr2, tpr2, threshold2 = metrics.roc_curve(y_tr, preds2)
roc_auc2 = metrics.auc(fpr2, tpr2)


probs1 = rf.predict_proba(Xbow_test_std)
preds1 = probs1[:,1]
fpr1, tpr1, threshold1 = metrics.roc_curve(y_test, preds1)
roc_auc1 = metrics.auc(fpr1, tpr1)
```

In [18]:
```
plt.title('Receiver Operating Characteristic on train and test')
plt.plot(fpr1, tpr1, 'b', label = 'AUC = %0.2f' % roc_auc1)
plt.plot(fpr2, tpr2, 'g', label = 'AUC = %0.2f' % roc_auc2)
#plt.plot(neighbors, auc1,'g')
#plt.plot(neighbors, auc2,'r')
plt.legend(loc = 'lower right')
plt.plot([0, 1], [0, 1],'r--')
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()
```



# UP-SAMPLING TECHNIQUE:

In [19]:
```python
#spilliting data

from sklearn.model_selection import train_test_split

# split the data set into train and test
X_1, X_test, y_1, y_test = train_test_split(df_tr_after_drop, yees, test_size=
0.3, random_state=42,stratify=yees)

# split the train data set into cross validation train and cross validation te
st
X_tr, X_cv, y_tr, y_cv = train_test_split(X_1, y_1, test_size=0.3)
```

In [20]:
```python
y_tr = pd.DataFrame(y_tr)
X_tr = pd.DataFrame(X_tr)
X_tr['target']=y_tr['target']
X_tr.head(2)
```

Out[20]:

| | ID | Attr1 | Attr2 | Attr3 | Attr4 | Attr5 | Attr6 | Attr7 | Attr8 | Att |
|---|---|---|---|---|---|---|---|---|---|---|
| **13653** | 13654 | -0.317430 | 0.76524 | 0.044714 | 1.1245 | -3.4901 | -1.903000 | -0.31743 | 0.30678 | 0.6549 |
| **23002** | 23003 | 0.076737 | 0.26259 | 0.448000 | 3.0413 | 54.0490 | 0.076737 | 0.09711 | 2.51650 | 1.0608 |

2 rows × 66 columns

In [21]:
```python
# Class count
X_tr = pd.DataFrame(X_tr)
print(type(X_tr))
count_class_0, count_class_1 = X_tr.target.value_counts()
print(count_class_0, count_class_1)

# Divide by class
df_class_0 = X_tr[X_tr.target == 0]
df_class_1 = X_tr[X_tr.target == 1]
print(count_class_0, count_class_1)


df_class_1_over = df_class_1.sample(count_class_0, replace=True)
X_tr = pd.concat([df_class_0, df_class_1_over], axis=0)

print('Random over-sampling:')
print(X_tr.target.value_counts())

X_tr.target.value_counts().plot(kind='bar', title='Count (target)');
```
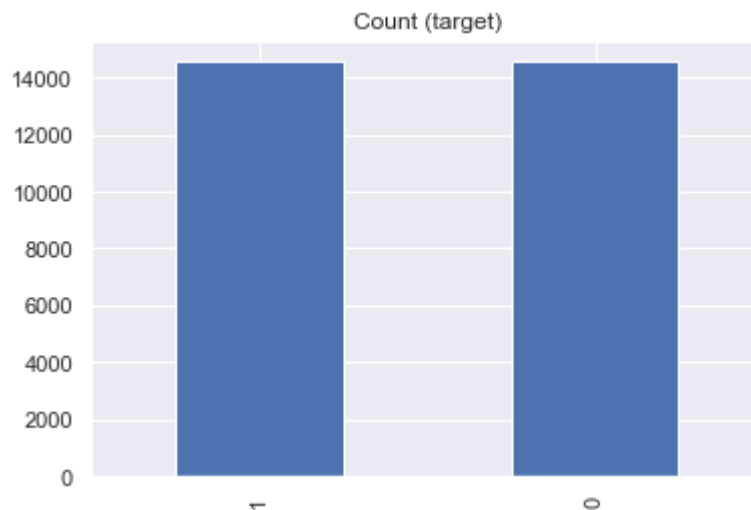
```
<class 'pandas.core.frame.DataFrame'>
14569 759
14569 759
Random over-sampling:
1    14569
0    14569
Name: target, dtype: int64
```



In [22]:
```python
y_tr=X_tr['target']
X_tr=X_tr.drop(['target'],axis=1)
```

```
In [23]: X_tr.fillna(X_tr.mean(), inplace=True)
         X_test.fillna(X_test.mean(), inplace=True)
         X_cv.fillna(X_cv.mean(), inplace=True)


         sc = StandardScaler(with_mean=True)
         Xbow_tr_std = sc.fit_transform(X_tr)
         Xbow_test_std = sc.transform(X_test)
         Xbow_cv_std = sc.transform(X_cv)
```

# HYPERPARAMETER TUNING WITH DEPTH AND NO.OF ESTIMATERS for XGBClassifier

```
In [24]: #code for hyperparameter tuning
         import numpy as np
         hyper1 = [5, 10, 50, 100, 200,400, 600, 800]
         hyper2 = [2, 4, 6, 8, 10, 12, 14, 16, 18]

         auc1=np.empty((8,9))
         auc2=np.empty((8,9))
         l=0
         for j in hyper1:
             m=0
             for k in hyper2:

                 model = XGBClassifier(n_estimators=j,max_depth=k)
                 model.fit(Xbow_tr_std, y_tr)

                 probs = model.predict_proba(Xbow_tr_std)
                 preds = probs[:,1]
                 roc_auc1=metrics.roc_auc_score(y_tr, preds)
                 auc1[l][m]=(roc_auc1)


                 probs = model.predict_proba(Xbow_cv_std)
                 preds = probs[:,1]
                 roc_auc2=metrics.roc_auc_score(y_cv, preds)
                 auc2[l][m]=(roc_auc2)
                 m=m+1

             l=l+1
```
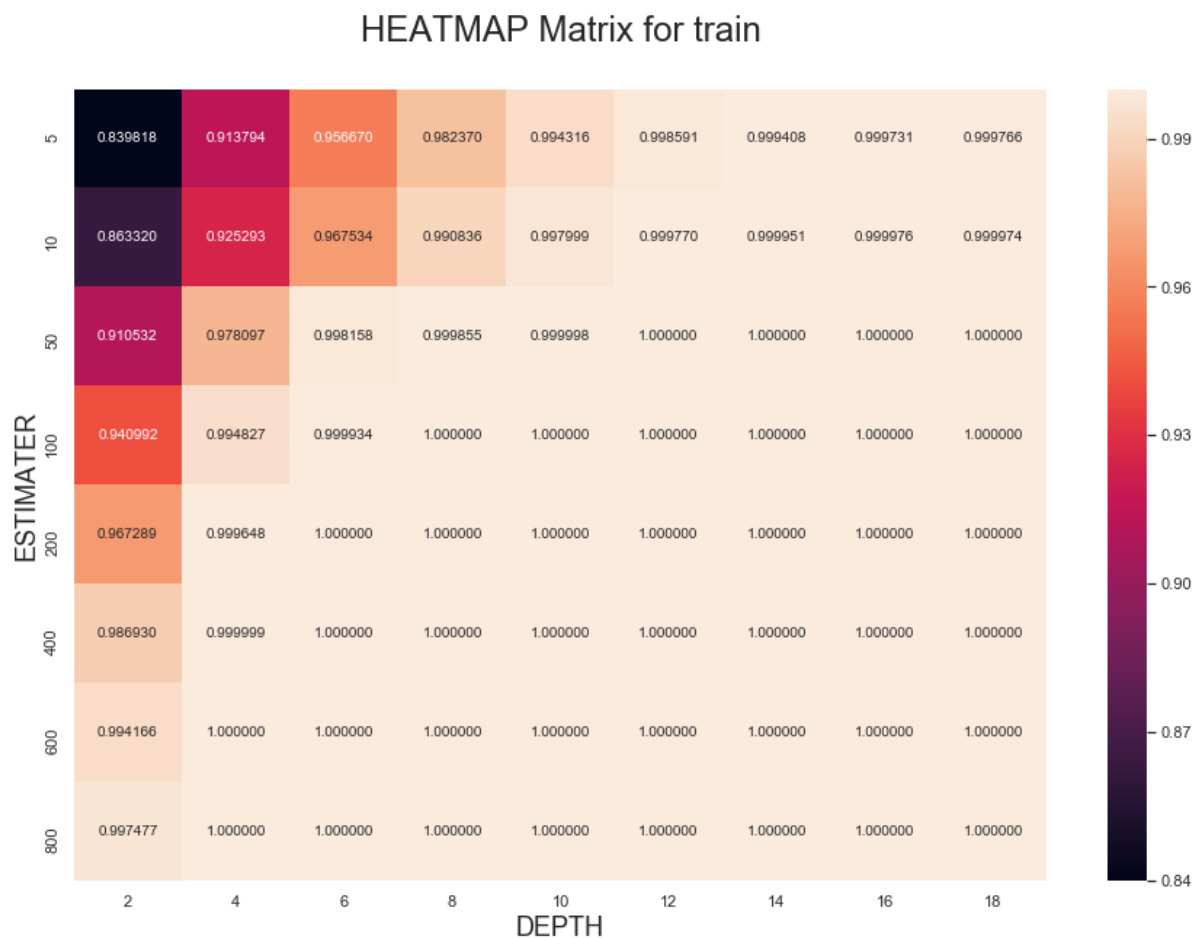
# PLOTTING SEABORN HEATMAP:

In [25]:
```python
# Code for drawing seaborn heatmaps

df_heatmap =pd.DataFrame(auc1, hyper1, hyper2 )
fig = plt.figure(figsize=(15,10))
ax = sns.heatmap(df_heatmap, annot=True, fmt="f")

plt.ylabel('ESTIMATER',size=18)
plt.xlabel('DEPTH',size=18)
plt.title("HEATMAP Matrix for train\n",size=24)
plt.show()
```

## HEATMAP Matrix for train

| ESTIMATER \ DEPTH | 2 | 4 | 6 | 8 | 10 | 12 | 14 | 16 | 18 |
|---|---|---|---|---|---|---|---|---|---|
| 5 | 0.839818 | 0.913794 | 0.956670 | 0.982370 | 0.994316 | 0.998591 | 0.999408 | 0.999731 | 0.999766 |
| 10 | 0.863320 | 0.925293 | 0.967534 | 0.990836 | 0.997999 | 0.999770 | 0.999951 | 0.999976 | 0.999974 |
| 50 | 0.910532 | 0.978097 | 0.998158 | 0.999855 | 0.999998 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 100 | 0.940992 | 0.994827 | 0.999934 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 200 | 0.967289 | 0.999648 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 400 | 0.986930 | 0.999999 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 600 | 0.994166 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 800 | 0.997477 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |

In [26]:
```python
# Code for drawing seaborn heatmaps

df_heatmap =pd.DataFrame(auc2, hyper1, hyper2 )
fig = plt.figure(figsize=(15,10))
ax = sns.heatmap(df_heatmap, annot=True, fmt="f")

plt.ylabel('ESTIMATER',size=18)
plt.xlabel('DEPTH',size=18)
plt.title("HEATMAP Matrix for CV\n",size=24)
plt.show()
```
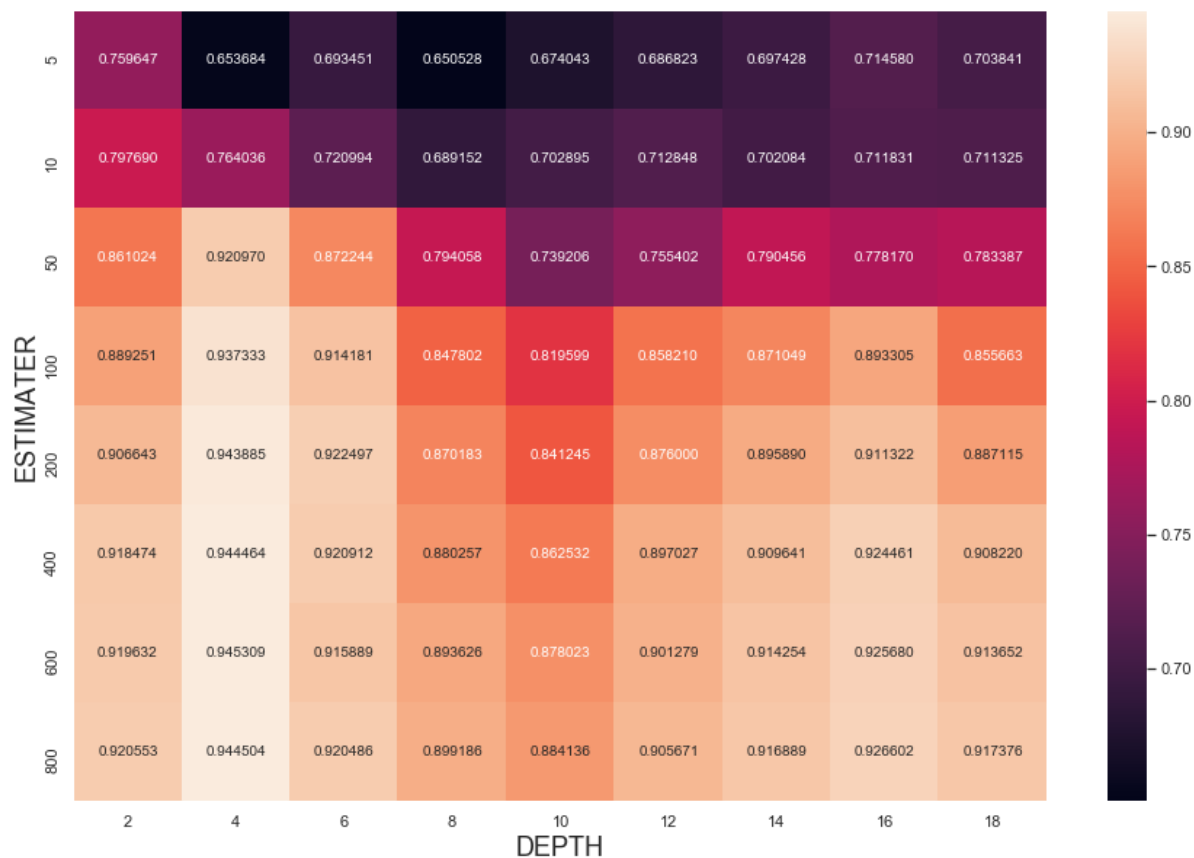
## HEATMAP Matrix for CV

| ESTIMATER \ DEPTH | 2 | 4 | 6 | 8 | 10 | 12 | 14 | 16 | 18 |
|---|---|---|---|---|---|---|---|---|---|
| 5 | 0.759647 | 0.653684 | 0.693451 | 0.650528 | 0.674043 | 0.686823 | 0.697428 | 0.714580 | 0.703841 |
| 10 | 0.797690 | 0.764036 | 0.720994 | 0.689152 | 0.702895 | 0.712848 | 0.702084 | 0.711831 | 0.711325 |
| 50 | 0.861024 | 0.920970 | 0.872244 | 0.794058 | 0.739206 | 0.755402 | 0.790456 | 0.778170 | 0.783387 |
| 100 | 0.889251 | 0.937333 | 0.914181 | 0.847802 | 0.819599 | 0.858210 | 0.871049 | 0.893305 | 0.855663 |
| 200 | 0.906643 | 0.943885 | 0.922497 | 0.870183 | 0.841245 | 0.876000 | 0.895890 | 0.911322 | 0.887115 |
| 400 | 0.918474 | 0.944464 | 0.920912 | 0.880257 | 0.862532 | 0.897027 | 0.909641 | 0.924461 | 0.908220 |
| 600 | 0.919632 | 0.945309 | 0.915889 | 0.893626 | 0.878023 | 0.901279 | 0.914254 | 0.925680 | 0.913652 |
| 800 | 0.920553 | 0.944504 | 0.920486 | 0.899186 | 0.884136 | 0.905671 | 0.916889 | 0.926602 | 0.917376 |

# FITTING AND TESTING MODEL ON OUR SPLITTED TEST DATA:

```
In [27]: rf = RandomForestClassifier(n_estimators=100,max_depth=6)

         # fitting the model
         rf.fit(Xbow_tr_std, y_tr)

         # predict the response
         pred = rf.predict(Xbow_test_std)

         # evaluate accuracy
         acc = accuracy_score(y_test, pred) * 100

         precision_score1=precision_score(y_test, pred )

         recall_score1=recall_score(y_test, pred )

         f1 = f1_score(y_test, pred)

         print('\nThe accuracy of the Random forest classifier for n_estimaters=%f and
          Depth = %f is %f%%' % (100,6, acc))



         print('\nThe precision_score of the  Random forest classifier  for n_estimater
         s=%d and Depth = %d is %f' % (100,6,precision_score1))

         print('\nThe recall_score of the  Random forest classifier  for n_estimaters=%
         d and Depth = %d is %f' % (100,6,recall_score1))

         print('\nThe f1_score of the  Random forest classifier  for n_estimaters=%d an
         d Depth = %d is %f' % (100,6,f1))
```

The accuracy of the Random forest classifier for n_estimaters=100.000000 and
Depth = 6.000000 is 85.572722%

The precision_score of the  Random forest classifier  for n_estimaters=100 an
d Depth = 6 is 0.210289

The recall_score of the  Random forest classifier  for n_estimaters=100 and D
epth = 6 is 0.721854

The f1_score of the  Random forest classifier  for n_estimaters=100 and Depth
= 6 is 0.325697

# PLOTTING CONFUSION MATRIX:

In [28]:
```python
# Code for drawing seaborn heatmaps
class_names = ['negative','positive']
df_heatmap = pd.DataFrame(confusion_matrix(y_test, pred), index=class_names, c
olumns=class_names )
fig = plt.figure(figsize=(10,7))
heatmap = sns.heatmap(df_heatmap, annot=True, fmt="d")

# Setting tick labels for heatmap
heatmap.yaxis.set_ticklabels(heatmap.yaxis.get_ticklabels(), rotation=0, ha='r
ight', fontsize=14)
heatmap.xaxis.set_ticklabels(heatmap.xaxis.get_ticklabels(), rotation=0, ha='r
ight', fontsize=14)
plt.ylabel('Predicted label',size=18)
plt.xlabel('True label',size=18)
plt.title("Confusion Matrix\n",size=24)
plt.show()
```
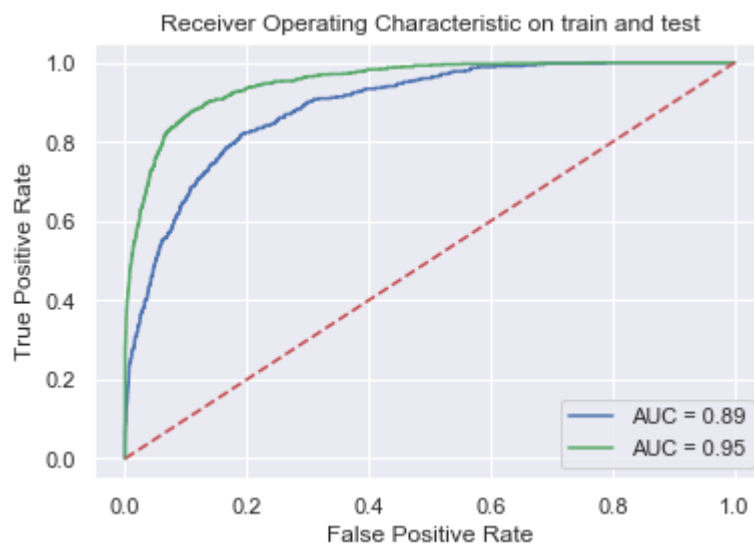


# PLOTTING AUC_ROC CURVE FOR TRAIN AND TEST DATA:

In [29]:
```python
rf.fit(Xbow_tr_std, y_tr)
probs2 = rf.predict_proba(Xbow_tr_std)
preds2 = probs2[:,1]
fpr2, tpr2, threshold2 = metrics.roc_curve(y_tr, preds2)
roc_auc2 = metrics.auc(fpr2, tpr2)


probs1 = rf.predict_proba(Xbow_test_std)
preds1 = probs1[:,1]
fpr1, tpr1, threshold1 = metrics.roc_curve(y_test, preds1)
roc_auc1 = metrics.auc(fpr1, tpr1)
```

In [30]:
```python
plt.title('Receiver Operating Characteristic on train and test')
plt.plot(fpr1, tpr1, 'b', label = 'AUC = %0.2f' % roc_auc1)
plt.plot(fpr2, tpr2, 'g', label = 'AUC = %0.2f' % roc_auc2)
#plt.plot(neighbors, auc1,'g')
#plt.plot(neighbors, auc2,'r')
plt.legend(loc = 'lower right')
plt.plot([0, 1], [0, 1],'r--')
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()
```



# HYPERPARAMETER TUNING WITH DEPTH AND NO.OF ESTIMATERS for XGBClassifier

```
In [31]:  #code for hyperparameter tuning
          import numpy as np
          hyper1 =  [5, 10, 50, 100, 200,400, 600, 800]
          hyper2 =  [2, 4, 6, 8, 10, 12, 14, 16, 18]

          auc1=np.empty((8,9))
          auc2=np.empty((8,9))
          l=0
          for j in hyper1:
              m=0
              for k in hyper2:

                  model = XGBClassifier(n_estimators=j,max_depth=k)
                  model.fit(Xbow_tr_std, y_tr)

                  probs = model.predict_proba(Xbow_tr_std)
                  preds = probs[:,1]
                  roc_auc1=metrics.roc_auc_score(y_tr, preds)
                  auc1[l][m]=(roc_auc1)


                  probs = model.predict_proba(Xbow_cv_std)
                  preds = probs[:,1]
                  roc_auc2=metrics.roc_auc_score(y_cv, preds)
                  auc2[l][m]=(roc_auc2)
                  m=m+1

              l=l+1
```
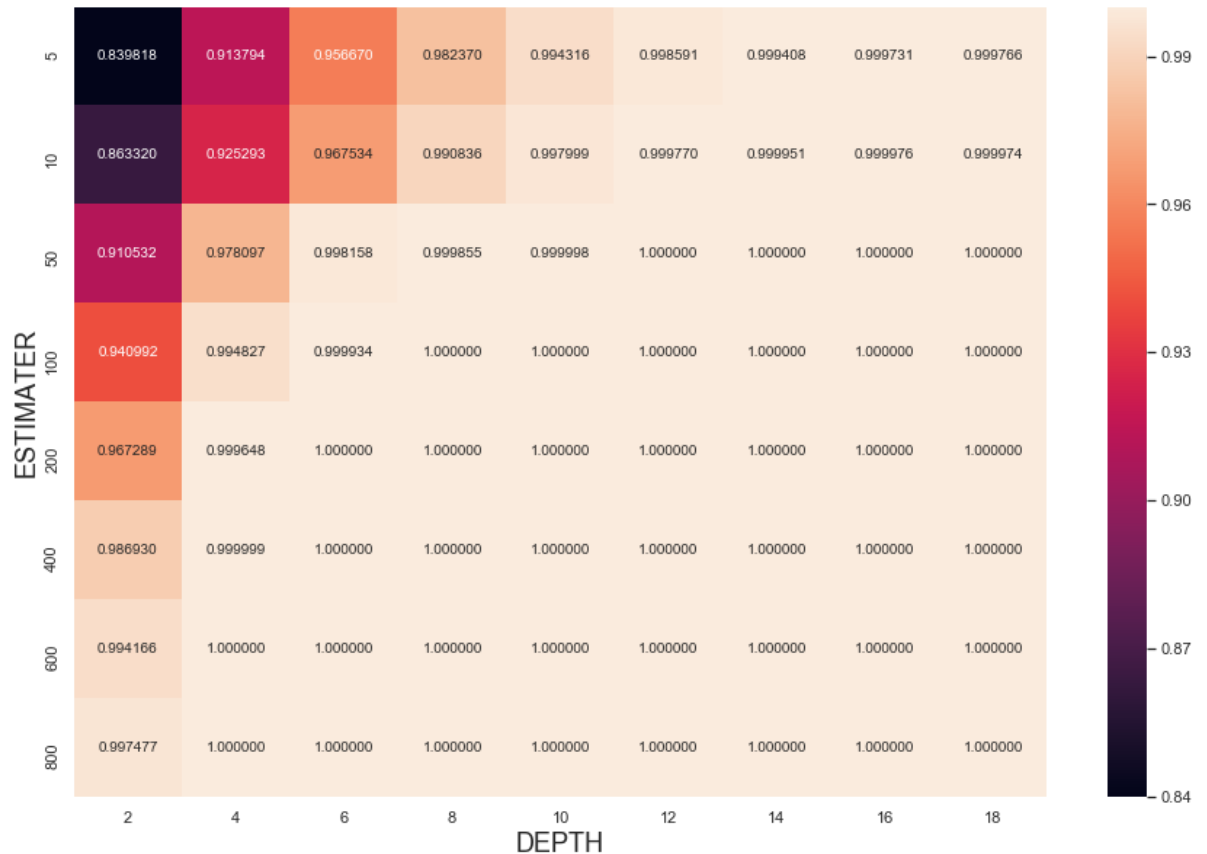
# PLOTTING SEABORN HEATMAP:

In [32]:
```python
# Code for drawing seaborn heatmaps

df_heatmap =pd.DataFrame(auc1, hyper1, hyper2 )
fig = plt.figure(figsize=(15,10))
ax = sns.heatmap(df_heatmap, annot=True, fmt="f")

plt.ylabel('ESTIMATER',size=18)
plt.xlabel('DEPTH',size=18)
plt.title("HEATMAP Matrix for train\n",size=24)
plt.show()
```
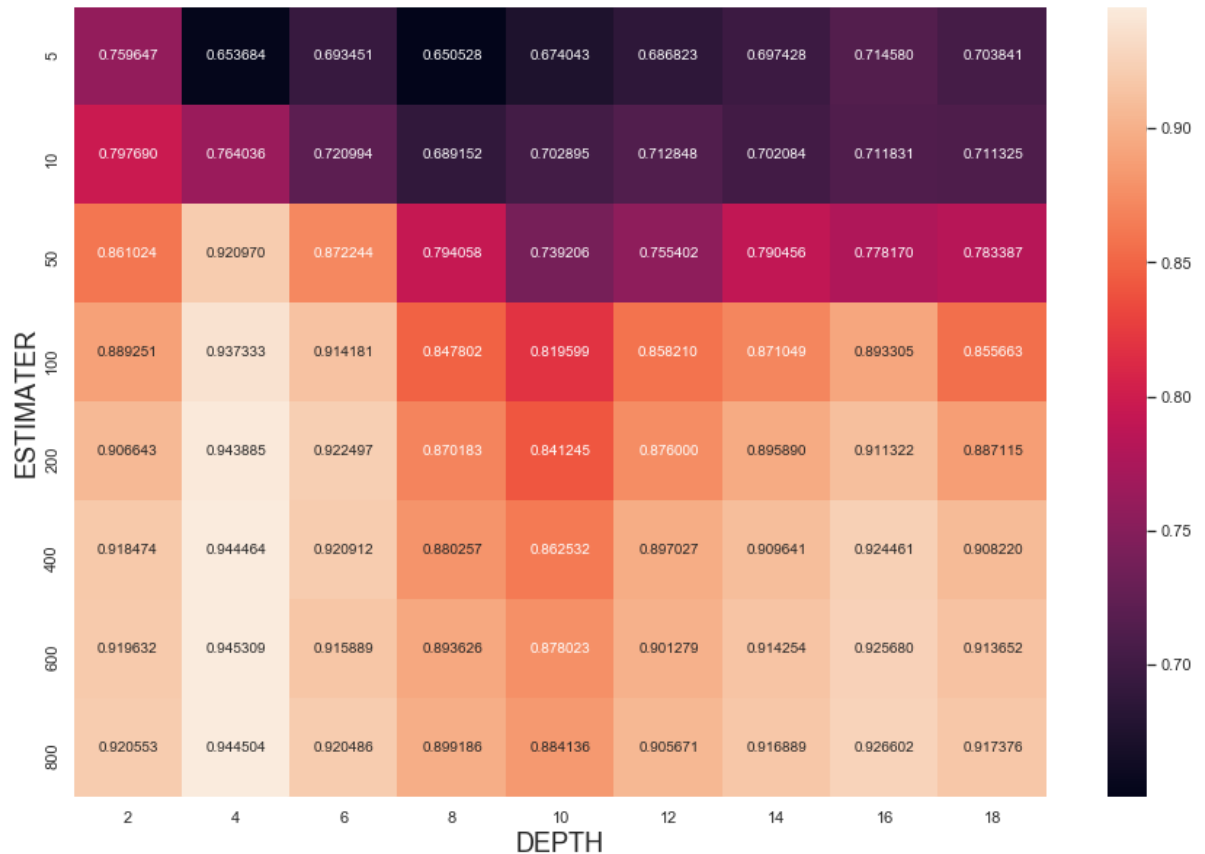
## HEATMAP Matrix for train

| ESTIMATER \ DEPTH | 2 | 4 | 6 | 8 | 10 | 12 | 14 | 16 | 18 |
|---|---|---|---|---|---|---|---|---|---|
| 5 | 0.839818 | 0.913794 | 0.956670 | 0.982370 | 0.994316 | 0.998591 | 0.999408 | 0.999731 | 0.999766 |
| 10 | 0.863320 | 0.925293 | 0.967534 | 0.990836 | 0.997999 | 0.999770 | 0.999951 | 0.999976 | 0.999974 |
| 50 | 0.910532 | 0.978097 | 0.998158 | 0.999855 | 0.999998 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 100 | 0.940992 | 0.994827 | 0.999934 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 200 | 0.967289 | 0.999648 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 400 | 0.986930 | 0.999999 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 600 | 0.994166 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| 800 | 0.997477 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 | 1.000000 |

```python
# Code for drawing seaborn heatmaps

df_heatmap =pd.DataFrame(auc2, hyper1, hyper2 )
fig = plt.figure(figsize=(15,10))
ax = sns.heatmap(df_heatmap, annot=True, fmt="f")

plt.ylabel('ESTIMATER',size=18)
plt.xlabel('DEPTH',size=18)
plt.title("HEATMAP Matrix for CV\n",size=24)
plt.show()
```

In [33]:

## HEATMAP Matrix for CV

| ESTIMATER \ DEPTH | 2 | 4 | 6 | 8 | 10 | 12 | 14 | 16 | 18 |
|---|---|---|---|---|---|---|---|---|---|
| 5 | 0.759647 | 0.653684 | 0.693451 | 0.650528 | 0.674043 | 0.686823 | 0.697428 | 0.714580 | 0.703841 |
| 10 | 0.797690 | 0.764036 | 0.720994 | 0.689152 | 0.702895 | 0.712848 | 0.702084 | 0.711831 | 0.711325 |
| 50 | 0.861024 | 0.920970 | 0.872244 | 0.794058 | 0.739206 | 0.755402 | 0.790456 | 0.778170 | 0.783387 |
| 100 | 0.889251 | 0.937333 | 0.914181 | 0.847802 | 0.819599 | 0.858210 | 0.871049 | 0.893305 | 0.855663 |
| 200 | 0.906643 | 0.943885 | 0.922497 | 0.870183 | 0.841245 | 0.876000 | 0.895890 | 0.911322 | 0.887115 |
| 400 | 0.918474 | 0.944464 | 0.920912 | 0.880257 | 0.862532 | 0.897027 | 0.909641 | 0.924461 | 0.908220 |
| 600 | 0.919632 | 0.945309 | 0.915889 | 0.893626 | 0.878023 | 0.901279 | 0.914254 | 0.925680 | 0.913652 |
| 800 | 0.920553 | 0.944504 | 0.920486 | 0.899186 | 0.884136 | 0.905671 | 0.916889 | 0.926602 | 0.917376 |

# FITTING AND TESTING MODEL ON OUR SPLITTED TEST DATA:

```
In [39]:  rf = RandomForestClassifier(n_estimators=100,max_depth=6)

          # fitting the model
          rf.fit(Xbow_tr_std, y_tr)

          # predict the response
          pred = rf.predict(Xbow_test_std)

          # evaluate accuracy
          acc = accuracy_score(y_test, pred) * 100

          precision_score1=precision_score(y_test, pred )

          recall_score1=recall_score(y_test, pred )

          f1 = f1_score(y_test, pred)

          print('\nThe accuracy of the Random forest classifier for n_estimaters=%f and
           Depth = %f is %f%%' % (100,6, acc))



          print('\nThe precision_score of the  Random forest classifier  for n_estimater
          s=%d and Depth = %d is %f' % (100,6,precision_score1))

          print('\nThe recall_score of the  Random forest classifier  for n_estimaters=%
          d and Depth = %d is %f' % (100,6,recall_score1))

          print('\nThe f1_score of the  Random forest classifier  for n_estimaters=%d an
          d Depth = %d is %f' % (100,6,f1))
```

The accuracy of the Random forest classifier for n_estimaters=100.000000 and
Depth = 6.000000 is 86.638253%

The precision_score of the  Random forest classifier  for n_estimaters=100 an
d Depth = 6 is 0.222453

The recall_score of the  Random forest classifier  for n_estimaters=100 and D
epth = 6 is 0.708609

The f1_score of the  Random forest classifier  for n_estimaters=100 and Depth
= 6 is 0.338608

# PLOTTING CONFUSION MATRIX:

In [35]:
```python
# Code for drawing seaborn heatmaps
class_names = ['negative','positive']
df_heatmap = pd.DataFrame(confusion_matrix(y_test, pred), index=class_names, c
olumns=class_names )
fig = plt.figure(figsize=(10,7))
heatmap = sns.heatmap(df_heatmap, annot=True, fmt="d")

# Setting tick labels for heatmap
heatmap.yaxis.set_ticklabels(heatmap.yaxis.get_ticklabels(), rotation=0, ha='r
ight', fontsize=14)
heatmap.xaxis.set_ticklabels(heatmap.xaxis.get_ticklabels(), rotation=0, ha='r
ight', fontsize=14)
plt.ylabel('Predicted label',size=18)
plt.xlabel('True label',size=18)
plt.title("Confusion Matrix\n",size=24)
plt.show()
```



# PLOTTING AUC_ROC CURVE FOR TRAIN AND TEST DATA:

In [36]:
```python
rf.fit(Xbow_tr_std, y_tr)
probs2 = rf.predict_proba(Xbow_tr_std)
preds2 = probs2[:,1]
fpr2, tpr2, threshold2 = metrics.roc_curve(y_tr, preds2)
roc_auc2 = metrics.auc(fpr2, tpr2)


probs1 = rf.predict_proba(Xbow_test_std)
preds1 = probs1[:,1]
fpr1, tpr1, threshold1 = metrics.roc_curve(y_test, preds1)
roc_auc1 = metrics.auc(fpr1, tpr1)
```

In [37]:
```python
plt.title('Receiver Operating Characteristic on train and test')
plt.plot(fpr1, tpr1, 'b', label = 'AUC = %0.2f' % roc_auc1)
plt.plot(fpr2, tpr2, 'g', label = 'AUC = %0.2f' % roc_auc2)
#plt.plot(neighbors, auc1,'g')
#plt.plot(neighbors, auc2,'r')
plt.legend(loc = 'lower right')
plt.plot([0, 1], [0, 1],'r--')
plt.ylabel('True Positive Rate')
plt.xlabel('False Positive Rate')
plt.show()
```



# PREETY_TABLE_OBSERVATION:

In [40]:
```python
from prettytable import PrettyTable

x = PrettyTable()

x.field_names = ["Technique", "Model","DEPTH","ESTIMATER", "Precession","recall","F1","ACCURACY"]

x.add_row(["No_Sampling", "Decision Tree",     8,    50,     1.000000,    0.022075,     0.043197,    95.279702])

x.add_row(["No_Sampling", "XGB_Regressor",     8,    50,     0.950000 ,   0.041943 ,    0.080338,   95.364944])

x.add_row(["SMOTE", "Decision Tree",          10,    50,     0.290323,    0.655629 ,    0.402439,   90.602025])

x.add_row(["SMOTE", "XGB_Regressor",           8,   100,     0.248366,    0.671082  ,   0.362552,   88.609483])

x.add_row(["UP_SAMPLING", "Decision Tree",     6,   100,     0.210289,    0.721854  ,   0.325697,   85.572722])

x.add_row(["UP_SAMPLING", "XGB_Regressor",     8,   100,     0.208071,    0.728477  ,   0.323688,  85.306340])


print(x)
```

```
+-------------+---------------+-------+-----------+------------+----------+----------+-----------+
|  Technique  |     Model     | DEPTH | ESTIMATER | Precession |  recall  |    F1    |  ACCURACY |
+-------------+---------------+-------+-----------+------------+----------+----------+-----------+
| No_Sampling | Decision Tree |   8   |     50    |     1.0    | 0.022075 | 0.043197 | 95.279702 |
| No_Sampling | XGB_Regressor |   8   |     50    |    0.95    | 0.041943 | 0.080338 | 95.364944 |
|    SMOTE    | Decision Tree |   10  |     50    |  0.290323  | 0.655629 | 0.402439 | 90.602025 |
|    SMOTE    | XGB_Regressor |   8   |    100    |  0.248366  | 0.671082 | 0.362552 | 88.609483 |
| UP_SAMPLING | Decision Tree |   6   |    100    |  0.210289  | 0.721854 | 0.325697 | 85.572722 |
| UP_SAMPLING | XGB_Regressor |   8   |    100    |  0.208071  | 0.728477 | 0.323688 |  85.30634 |
+-------------+---------------+-------+-----------+------------+----------+----------+-----------+
```

# -------------------SAVING_RESULT_WITHOUT_SAMPLING---------------------

```
In [170]: df_tr1=pd.read_csv('train.csv')
          df_te1=pd.read_csv('test.csv')
          yees=df_tr['target']
          df_tr_after_drop=df_tr.drop(['target'],axis=1)
```

```
In [171]: df_tr_after_drop.fillna(df_tr_after_drop.mean(), inplace=True)
          df_te1.fillna(df_te1.mean(), inplace=True)


          Xbow_tr_std = sc.fit_transform(df_tr_after_drop)
          Xbow_test_std = sc.transform(df_te1)
```

```
In [173]: rf = XGBClassifier(n_estimators=150,max_depth=10)

          # fitting the model
          rf.fit(Xbow_tr_std, yees)

          # predict the response
          pred = rf.predict(Xbow_test_std)
```

```
In [174]: df_te['TARGET'] = pred
          df_te.to_csv('test_with_target_finel1.csv', index=True)
```

# -----------------------------SAVING_RESULT_WITH_SMOTE-------------------------------

```
In [175]: X_tr, y_tr = sm.fit_sample(df_tr_after_drop, yees)
          Xbow_tr_std = sc.fit_transform(X_tr)
```

```
In [176]: rf = XGBClassifier(n_estimators=150,max_depth=10)

          # fitting the model
          rf.fit(Xbow_tr_std, y_tr)

          # predict the response
          pred = rf.predict(Xbow_test_std)
```

```
In [177]: df_te['TARGET'] = pred
          df_te.to_csv('test_with_target_finel2.csv', index=True)
```

# ------SAVING_RESULT_WITH_OVERSAMPLING_MINORITY_CL--------

In [178]:
```python
# Class count
X_tr = pd.DataFrame(df_tr1)
print(type(X_tr))
count_class_0, count_class_1 = X_tr.target.value_counts()
print(count_class_0, count_class_1)

# Divide by class
df_class_0 = X_tr[X_tr.target == 0]
df_class_1 = X_tr[X_tr.target == 1]
print(count_class_0, count_class_1)


df_class_1_over = df_class_1.sample(count_class_0, replace=True)
X_tr = pd.concat([df_class_0, df_class_1_over], axis=0)

print('Random over-sampling:')
print(X_tr.target.value_counts())

X_tr.target.value_counts().plot(kind='bar', title='Count (target)');
```
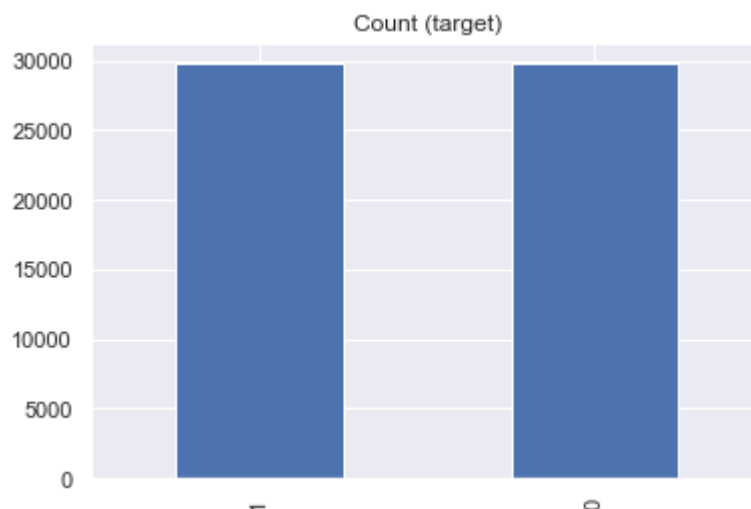
```
<class 'pandas.core.frame.DataFrame'>
29772 1511
29772 1511
Random over-sampling:
1    29772
0    29772
Name: target, dtype: int64
```


Count (target)

In [179]:
```python
y_tr=X_tr['target']
X_tr=X_tr.drop(['target'],axis=1)
```

In [180]:
```python
X_tr.fillna(X_tr.mean(), inplace=True)

Xbow_tr_std = sc.fit_transform(X_tr)
```

In [181]:
```python
rf = XGBClassifier(n_estimators=100,max_depth=8)

# fitting the model
rf.fit(Xbow_tr_std, y_tr)

# predict the response
pred = rf.predict(Xbow_test_std)
```

In [182]:
```python
df_te['TARGET'] = pred
df_te.to_csv('test_with_target_finel3.csv', index=True)
```

In [ ]: