

Disclaimer and acknowledgement

Author: H. NIU

July 21, 2021

The notebook was partly borrowed, adapted and function-added for needs from https://github.com/IntelliHQ/CardiacArrestMortality_ANZICS (https://github.com/IntelliHQ/CardiacArrestMortality_ANZICS), which holds the copyright and credits on that code. This is acknowledged the LICENSE with that repo is not copied along. I greatly appreciate the authors of that code for their efforts.

In [1]: 1 !python -V

Python 3.7.3

```
In [2]: 1 import pickle
2 import random
3 import pandas as pd
4 import numpy as np
5 import scipy.stats as stats
6 import sklearn
7 import matplotlib
8
9 from sklearn.base import BaseEstimator, TransformerMixin, ClassifierMixin
10 from sklearn.preprocessing import StandardScaler
11 from sklearn.pipeline import Pipeline, FeatureUnion
12
13 from sklearn.model_selection import train_test_split
14
15 from sklearn.metrics import accuracy_score, confusion_matrix, roc_curve, auc
16 #from sklearn.metrics.scorer import make_scorer
17
18 from lime.lime_tabular import LimeTabularExplainer
19 import matplotlib.pyplot as plt
20
21 pd.set_option('display.max_columns', None)
22
23 print('Numpy version:', np.__version__)
24 print('Pandas version:', pd.__version__)
25 print('sklearn version:', sklearn.__version__)
26 print('matplotlib version:', matplotlib.__version__)
```

Numpy version: 1.19.5
Pandas version: 0.25.1
sklearn version: 1.0
matplotlib version: 3.4.3

```
In [3]: 1 #Imports
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import itertools
6
7 #scikit-learn package (https://pypi.org/project/scikit-learn)
8 from sklearn.model_selection import train_test_split
9 from sklearn.ensemble import RandomForestClassifier
10 from sklearn.metrics import accuracy_score, confusion_matrix, roc_curve, auc
11 #from sklearn.metrics.scorer import make_scorer
12
13 #eli5 package (https://eli5.readthedocs.io/en/latest)
14 import eli5
15 from eli5.sklearn import PermutationImportance
16
17 #lime package (https://github.com/marcotcr/Lime)
18 import lime
19 import lime.lime_tabular
20
21 #shap package (https://github.com/slundberg/shap)
22 import shap
```

```
In [4]: 1 # %matplotlib inline
2 plt.style.use('ggplot')
3 plt.rcParams['figure.figsize'] = 12, 8
4 plt.rcParams['xtick.labelsize'] = 12
5 plt.rcParams['ytick.labelsize'] = 12
6 plt.rcParams['axes.titlesize'] = 18
```

Making machine learning pipeline

LIME explainers take in the original raw data for model explanation (for human understandable explanations), so the target model has to include the preprocessing steps. So the following pipeline is largely a repeat of what's done in the preprocessing and omitting the model fitting steps.

```
In [5]: 1 df = pd.read_excel(r'C:\Users\neo\Huijie\Code\Predicting CA Mortality\Data\1
2
3 df.head()
```

Out[5]:

	Gender	Age_value	BMI_value	Comorbidity_diabetes	Comorbidity_hypertension	Comorbidity_CA
0	Female	46.0	23.738662	0	0	
1	Male	72.0	24.341758	0	0	
2	Male	56.0	24.034610	0	0	
3	Female	64.0	29.048656	0	0	
4	Male	68.0	28.081633	0	0	

In [6]:

```

1 df['Age']='<12'
2 #df['Age']='>=65'
3 df.loc[(df['Age_value'] >= 12) & (df['Age_value'] < 40), 'Age'] = '12~40'
4 #df.loc[(df['Age_value'] < 40), 'Age'] = '<40'
5 df.loc[(df['Age_value'] >= 40) & (df['Age_value'] < 65), 'Age'] = '40~65'
6 df.loc[(df['Age_value'] >= 65), 'Age'] = '>=65'
7 #df.loc[(df['Age_value'] <=5), 'Age'] = '<=5'
8
9 df['BMI']='Underweight'
10 df.loc[(df['BMI_value'] >= 19) & (df['BMI_value'] < 24), 'BMI'] = 'Ideal'
11 df.loc[(df['BMI_value'] >= 24) & (df['BMI_value'] < 28), 'BMI'] = 'Overweigh'
12 df.loc[(df['BMI_value'] >= 28), 'BMI'] = 'Obese'
13
14 #df['ASA']='I-III'
15 #df.loc[(df['ASA_PS'] >= 4), 'ASA'] = 'IV-V'
16 df['ASA']='I'
17 df.loc[(df['ASA_PS'] >= 2) & (df['ASA_PS'] < 3), 'ASA'] = 'II'
18 df.loc[(df['ASA_PS'] >= 3) & (df['ASA_PS'] < 4), 'ASA'] = 'III'
19 df.loc[(df['ASA_PS'] >= 4) & (df['ASA_PS'] < 5), 'ASA'] = 'IV'
20 df.loc[(df['ASA_PS'] >= 5), 'ASA'] = 'V'
21
22 #df['Hemorrhage']='0'
23 df['Hemorrhage'] = '<200'
24 df.loc[(df['Hemorrhage_ml'] >= 200) & (df['Hemorrhage_ml'] < 800), 'Hemorrha'
25 df.loc[(df['Hemorrhage_ml'] >= 800), 'Hemorrhage'] = '>=800'
26
27 #df['Blood_transfusion']='0'
28 df['Blood_transfusion'] = '<200'
29 df.loc[(df['Blood_transfusion_ml'] >= 200) & (df['Blood_transfusion_ml'] < 8
30 df.loc[(df['Blood_transfusion_ml'] >= 800), 'Blood_transfusion'] = '>=800'
31
32 df['Epinephrine']='0'
33 df.loc[(df['Epinephrine_mg'] > 0) & (df['Epinephrine_mg'] <= 5), 'Epinephrin
34 df.loc[(df['Epinephrine_mg'] > 5), 'Epinephrine'] = '>5'
35
36 df['Atropine']='0'
37 df.loc[(df['Atropine_mg'] > 0) & (df['Atropine_mg'] <= 0.65), 'Atropine'] =
38 df.loc[(df['Atropine_mg'] > 0.65), 'Atropine'] = '>0.65'
39
40 df['Amiodarone']='0'
41 #df.loc[(df['Amiodarone_g'] > 0) & (df['Amiodarone_g'] <= 0.11), 'Amiodarone
42 df.loc[(df['Amiodarone_g'] > 0), 'Amiodarone'] = '0.04~0.3'
43
44 df['Ephedrine']='0'
45 #df.loc[(df['Ephedrine_mg'] > 0) & (df['Ephedrine_mg'] <= 6), 'Ephedrine'] =
46 df.loc[(df['Ephedrine_mg'] > 0), 'Ephedrine'] = '6~15'
47
48 df['Methoxamine']='0'
49 #df.loc[(df['Methoxamine_mg'] > 0) & (df['Methoxamine_mg'] <= 3), 'Methoxami
50 df.loc[(df['Methoxamine_mg'] > 0), 'Methoxamine'] = '1~35'
51
52 #df['CPR']='0'
53 #df.loc[(df['CPR_min'] > 0) & (df['CPR_min'] < 30), 'CPR'] = '<=30'
54 df.loc[(df['CPR_min'] < 30), 'CPR'] = '<30'
55 df.loc[(df['CPR_min'] >= 30) & (df['CPR_min'] < 60), 'CPR'] = '30~60'
56 df.loc[(df['CPR_min'] >= 60), 'CPR'] = '>=60'

```

```

57
58
59 df2=df[["Gender","Age","BMI","Comorbidity_diabetes","Comorbidity_hypertensio
60 "Comorbidity_pulmonary","Comorbidity_hepatic","Comorbidity_renal","Comorbidi
61 "Comorbidity_tumor","Surgical_type","Emergency","Trauma","Anaesthetic_type",
62 "ASA","Timing_arrest","Defibrillate","Cause_arrest","Hemorrhage","Blood_tran
63 "Epinephrine","Atropine","Amiodarone","Ephedrine","Methoxamine","CPR","Died"
64
65 df2.head()
66 df2.info()

```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 29 columns):
Gender                150 non-null object
Age                  150 non-null object
BMI                  150 non-null object
Comorbidity_diabetes  150 non-null int64
Comorbidity_hypertension  150 non-null int64
Comorbidity_CAD      150 non-null int64
Comorbidity_pulmonary  150 non-null int64
Comorbidity_hepatic   150 non-null int64
Comorbidity_renal     150 non-null int64
Comorbidity_neurological  150 non-null int64
Comorbidity_tumor     150 non-null int64
Surgical_type        150 non-null object
Emergency            150 non-null int64
Trauma               150 non-null int64
Anaesthetic_type     150 non-null object
Operative_position   150 non-null object
ASA                  150 non-null object
Timing_arrest        150 non-null object
Defibrillate         150 non-null int64
Cause_arrest         150 non-null object
Hemorrhage           150 non-null object
Blood_transfusion    150 non-null object
Epinephrine          150 non-null object
Atropine              150 non-null object
Amiodarone           150 non-null object
Ephedrine            150 non-null object
Methoxamine          150 non-null object
CPR                  150 non-null object
Died                 150 non-null int64
dtypes: int64(12), object(17)
memory usage: 34.1+ KB

```

In [7]:

```

1 df=df2
2 X=df
3 X['Gender'] = X.Gender.map({'Male':1, 'Female':0}) # M -> 1, F -> 0
4 X['Age'] = X.Age.map({'<12':0, '12~40':1, '40~65':2, '>=65':3}) #
5 X['BMI'] = X.BMI.map({'Underweight':0, 'Ideal':1, 'Overweight':2, 'Obese':3})
6 X['ASA'] = X.ASA.map({'I':0, 'II':1, 'III':2, 'IV':3, 'V':4}) #
7 X['Hemorrhage'] = X.Hemorrhage.map({'<200':0, '200~800':1, '>=800':2}) #
8 X['Blood_transfusion'] = X.Blood_transfusion.map({'<200':0, '200~800':1, '>=800':2}) #
9 X['Epinephrine'] = X.Epinephrine.map({'0':0, '<=5':1, '>5':2}) #
10 X['Atropine'] = X.Atropine.map({'0':0, '<=0.65':1, '>0.65':2}) #
11 X['Amiodarone'] = X.Amiodarone.map({'0':0, '0.04~0.3':1}) #
12 X['Ephedrine'] = X.Ephedrine.map({'0':0, '6~15':1}) #
13 X['Methoxamine'] = X.Methoxamine.map({'0':0, '1~35':1}) #
14 X['CPR'] = X.CPR.map({'<30':0, '30~60':1, '>=60':2}) #
15 X['Surgical_type'] = X.Surgical_type.map({'Abdominal':0, 'Neurosurgery':1, 'Thoracic':2, 'Throat':3, 'Others':4}) #
16 # Abdominal->0, Neurosurgery->1, Thoracic->2, Throat->3, Others->4
17 X['Anaesthetic_type'] = X.Anaesthetic_type.map({'General':1, 'Local':0}) #
18 X['Operative_position'] = X.Operative_position.map({'Left lateral decubitus':1, 'Lithotomy':2, 'Supine':3}) #
19 # Left lateral decubitus->1, Lithotomy->2, Supine->3
20 X['Timing_arrest'] = X.Timing_arrest.map({'Induction':0, 'Intubation':1, 'Surveillance':2}) #
21 X['Cause_arrest'] = X.Cause_arrest.map({'Anesthesia':0, 'Comorbidity':1, 'Surveillance':2}) #
22 df=X

```

In [8]:

```

1 def split_train_test(data, test_ratio):
2     shuffled_indices = np.random.permutation(len(data))
3     test_set_size = int(len(data) * test_ratio)
4     test_indices = shuffled_indices[:test_set_size]
5     train_indices = shuffled_indices[test_set_size:]
6     return data.iloc[train_indices], data.iloc[test_indices]

```

In [9]:

```

1 #Now split the data into training and test set(85-15 split)
2 data_train, data_test = train_test_split(df, test_size=.25,
3                                           stratify=df.Gender, random_state=0)

```

In [10]:

```
1 data_train.head()
```

Out[10]:

	Gender	Age	BMI	Comorbidity_diabetes	Comorbidity_hypertension	Comorbidity_CAD	Comorbidity
85	1	2	1	0	1	0	
105	0	2	3	0	0	0	
65	0	1	1	0	0	1	
142	1	2	3	0	1	1	
149	1	2	2	0	1	1	

Data preprocessing

```
In [11]: 1 categorical_vars = ["Gender", "Age", "BMI", "Comorbidity_diabetes", "Comorbidity_
2 "Comorbidity_pulmonary", "Comorbidity_hepatic", "Comorbidity_renal", "Comorbidi
3 "Comorbidity_tumor", "Surgical_type", "Emergency", "Trauma", "Anaesthetic_type",
4 "ASA", "Timing_arrest", "Defibrillate", "Cause_arrest", "Hemorrhage", "Blood_tran
5 "Epinephrine", "Atropine", "Amiodarone", "Ephedrine", "Methoxamine", "CPR"]
6
7 categorical_var_idx = [idx for idx, _ in enumerate(list(data_train.columns))
8 print(categorical_var_idx)]
```

```
[0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21,
22, 23, 24, 25, 26, 27]
```

```
In [12]: 1 categorical_vars2 = ["Gender", "Age", "BMI", "Comorbidity_diabetes", "Comorbidity_
2 "Comorbidity_pulmonary", "Comorbidity_hepatic", "Comorbidity_renal", "Comorbidi
3 "Comorbidity_tumor", "Surgical_type", "Emergency", "Trauma", "Anaesthetic_type",
4 "ASA", "Timing_arrest", "Defibrillate", "Cause_arrest", "Hemorrhage", "Blood_tran
5 "Epinephrine", "Atropine", "Amiodarone", "Ephedrine", "Methoxamine", "CPR", "Died"
6
7 categorical_var_idx2 = [idx for idx, _ in enumerate(list(data_train.columns))
8 print(categorical_var_idx2)]
```

```
[0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21,
22, 23, 24, 25, 26, 27, 28]
```

```
In [13]: 1 # to make a custom transformer to fit into a pipeline
2 class Vars_selector(BaseEstimator, TransformerMixin):
3     '''Return a subset of variables in a numpy array based on indecies'''
4     def __init__(self, var_idx):
5         '''var_idx is a list of categorical variables indecies'''
6         self.var_idx = var_idx
7
8     def fit(self, X, y=None):
9         return self
10
11    def transform(self, X):
12        '''returns a dataframe with selected variables'''
13        return np.array(X[:, self.var_idx])
```

```
In [14]: 1 class Cat_vars_encoder(BaseEstimator, TransformerMixin):
2     '''Return the transformed categorical variables based on indecies'''
3     def fit(self, X, y=None):
4         return self
5
6     def transform(self, X):
7         # arf column index
8         # arf_idx = categorical_vars.index('arf')
9         # X[:, arf_idx] = np.array(pd.Series(X[:, arf_idx]).map({1:1, 2:0}))
10        return X
```

Transform data in a pipeline

```
In [15]: 1 # categorical variables preprocessing
2 cat_vars_pipeline = Pipeline([
3     ('selector', Vars_selector(categorical_var_idxes)),
4     ('encoder', Cat_vars_encoder())
5 ])
```

```
In [16]: 1 # categorical variables preprocessing
2 cat_vars_pipeline2 = Pipeline([
3     ('selector', Vars_selector(categorical_var_idxes2)),
4     ('encoder', Cat_vars_encoder())
5 ])
```

```
In [17]: 1 continuous_vars = [] #['Age', 'BMI', 'Hemorrhage_ml', 'Blood_transfusion', '
2         #'Amiodarone_g', 'Ephedrine_mg', 'Methoxamine_mg', 'CPR_min']
```

To transform the variables in one step.

```
In [18]: 1 preproc_pipeline = FeatureUnion(transformer_list=[
2     ('cat_pipeline', cat_vars_pipeline)
3     #, ('cont_pipeline', cont_vars_pipeline)
4 ])
5
6
7 preproc_pipeline2 = FeatureUnion(transformer_list=[
8     ('cat_pipeline', cat_vars_pipeline2)
9     #, ('cont_pipeline', cont_vars_pipeline)
10 ])
```

```
In [19]: 1 data_train_X = pd.DataFrame(preproc_pipeline.fit_transform(data_train),
2                             columns=categorical_vars + continuous_vars)
3
4 data_train_X2 = pd.DataFrame(preproc_pipeline2.fit_transform(data_train),
5                             columns=categorical_vars2 + continuous_vars)
```

```
In [20]: 1 data_test_X = pd.DataFrame(preproc_pipeline2.fit_transform(data_test),
2                             columns=categorical_vars2 + continuous_vars)
3
4 data_test_X.head()
```

Out[20]:

	Gender	Age	BMI	Comorbidity_diabetes	Comorbidity_hypertension	Comorbidity_CAD	Comorbi
0	0	2	0	0	0	1	
1	0	1	1	0	0	0	
2	0	2	2	0	0	0	
3	0	1	2	0	1	1	
4	1	2	2	0	1	1	


```
In [21]: 1 print(data_train_X2.shape)
```

```
(112, 29)
```

Making explanation

1. A feature importance explainer

<https://machinelearningmastery.com/calculate-feature-importance-with-python/>
(<https://machinelearningmastery.com/calculate-feature-importance-with-python/>)

To restore the models and fit them into pipeline.

```
In [22]: 1 with open(r'C:\Users\neo\Huijie\Code\Predicting CA Mortality\Models\rf_clf_f  
2         rf_clf = pickle.load(f)  
3         model=rf_clf
```

Trying to unpickle estimator DecisionTreeClassifier from version 0.24.0 when using version 1.0. This might lead to breaking code or invalid results. Use at your own risk. For more info please refer to:

https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations (https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations)

Trying to unpickle estimator RandomForestClassifier from version 0.24.0 when using version 1.0. This might lead to breaking code or invalid results. Use at your own risk. For more info please refer to:

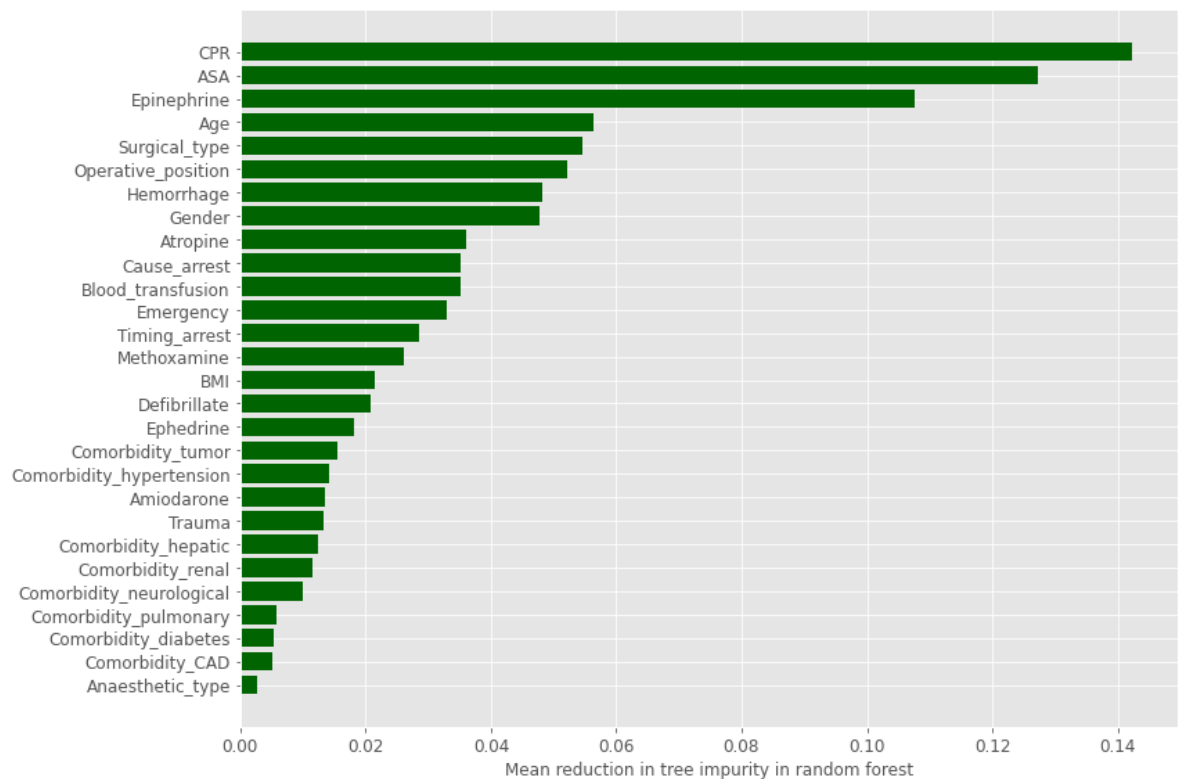
https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations (https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations)

```
In [23]: 1 #data=data_train.drop('Died',1)  
2         data=data_train_X
```

```

In [24]: 1 import matplotlib.patches as patches
2 import matplotlib.colors as colors
3 import math
4
5
6 # Feature importance dataframe
7 imp_df = pd.DataFrame({'feature': data.columns.values,
8                        'importance': model.feature_importances_})
9
10 # Reorder by importance
11 ordered_df = imp_df.sort_values(by='importance')
12 imp_range=range(1,len(imp_df.index)+1)
13
14 ## Barplot with confidence intervals
15 height = ordered_df['importance']
16 bars = ordered_df['feature']
17 y_pos = np.arange(len(bars))
18
19 # Create horizontal bars
20 plt.barh(y_pos, height, color='darkgreen')
21
22 # Create names on the y-axis
23 plt.yticks(y_pos, bars)
24
25 plt.xlabel("Mean reduction in tree impurity in random forest")
26
27 plt.tight_layout()
28
29
30 # Show graphic
31 plt.show()

```



In [25]:

```

1  # Feature importance based on TRAINING set
2
3  perm_test = PermutationImportance(model, scoring=make_scorer(roc_auc_score),
4                                     n_iter=50, random_state=0, cv="prefit")
5
6  # fit and see the permuation importances
7  perm_test.fit(data, data_train['Died'])
8
9  imp_df = eli5.explain_weights_df(perm_test)
10 label_df = pd.DataFrame({'feature': [ "x" + str(i) for i in range(len(data.c
11 imp_df = pd.merge(label_df, imp_df, on='feature', how='inner', validate="one
12
13 # Reorder by importance
14 ordered_df = imp_df.sort_values(by='weight')
15 imp_range=range(1,len(imp_df.index)+1)
16
17
18 ## Barplot with confidence intervals
19
20 height = ordered_df['weight']
21 bars = ordered_df['feature_name']
22 ci = 1.96 * ordered_df['std']
23 y_pos = np.arange(len(bars))
24
25 # Create horizontal bars
26 plt.barh(y_pos, height, xerr=ci, color='darkgreen')
27
28 # Create names on the y-axis
29 plt.yticks(y_pos, bars)
30
31 plt.xlabel("Permutation feature importance training set (decrease in AUC)")
32 plt.tight_layout()
33
34 # Show graphic
35 plt.show()
36

```

X has feature names, but RandomForestClassifier was fitted without feature names

X has feature names, but RandomForestClassifier was fitted without feature names

X has feature names, but RandomForestClassifier was fitted without feature names

X has feature names, but RandomForestClassifier was fitted without feature names

X has feature names, but RandomForestClassifier was fitted without feature names

X has feature names, but RandomForestClassifier was fitted without feature names

X has feature names, but RandomForestClassifier was fitted without feature names

X has feature names, but RandomForestClassifier was fitted without feature names

X has feature names, but RandomForestClassifier was fitted without feature names

X has feature names, but RandomForestClassifier was fitted without feature names

X has feature names, but RandomForestClassifier was fitted without feature na

2. Preparing a LIME explainer

To restore the models and fit them into pipeline.

```
In [26]: 1 with open(r'C:\Users\neo\Huijie\Code\Predicting CA Mortality\Models\ensemble_
2         ensemble_clf = pickle.load(f)
```

Trying to unpickle estimator LogisticRegression from version 0.24.0 when using version 1.0. This might lead to breaking code or invalid results. Use at your own risk. For more info please refer to:

https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations (https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations)

Trying to unpickle estimator DecisionTreeClassifier from version 0.24.0 when using version 1.0. This might lead to breaking code or invalid results. Use at your own risk. For more info please refer to:

https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations (https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations)

Trying to unpickle estimator RandomForestClassifier from version 0.24.0 when using version 1.0. This might lead to breaking code or invalid results. Use at your own risk. For more info please refer to:

https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations (https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations)

Trying to unpickle estimator DecisionTreeRegressor from version 0.24.0 when using version 1.0. This might lead to breaking code or invalid results. Use at your own risk. For more info please refer to:

https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations (https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations)

Trying to unpickle estimator DummyClassifier from version 0.24.0 when using version 1.0. This might lead to breaking code or invalid results. Use at your own risk. For more info please refer to:

https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations (https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations)

Trying to unpickle estimator GradientBoostingClassifier from version 0.24.0 when using version 1.0. This might lead to breaking code or invalid results. Use at your own risk. For more info please refer to:

https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations (https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations)

Trying to unpickle estimator SVC from version 0.24.0 when using version 1.0.

This might lead to breaking code or invalid results. Use at your own risk. For more info please refer to:

https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations (https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations)

Trying to unpickle estimator AdaBoostClassifier from version 0.24.0 when using version 1.0. This might lead to breaking code or invalid results. Use at your own risk. For more info please refer to:

https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations (https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations)

Trying to unpickle estimator LabelEncoder from version 0.24.0 when using version 1.0. This might lead to breaking code or invalid results. Use at your own risk. For more info please refer to:

https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations (https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations)

[ence.html#security-maintainability-limitations](#))

Trying to unpickle estimator VotingClassifier from version 0.24.0 when using version 1.0. This might lead to breaking code or invalid results. Use at your own risk. For more info please refer to:

https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations (https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations)

```
In [27]: 1  ## for ensemble_clf:
2  class Ensemble_clf(BaseEstimator, ClassifierMixin):
3      def fit(self, X, y=None):
4          return ensemble_clf
5      def predict_proba(self, X):
6          return ensemble_clf.predict_proba(X)
```

```
In [28]: 1  ## for ensemble_clf
2  ensemble_pipeline = Pipeline([
3      ('preprocessing', FeatureUnion(transformer_list=[
4          ('cat_pipeline', cat_vars_pipeline)
5          ])),
6      ('ensemble_clf', Ensemble_clf())
7  ])
```

```
In [29]: 1 ensemble_pipeline.fit(data_train)
```

```
Out[29]: Pipeline(steps=[('preprocessing',  
                           FeatureUnion(transformer_list=[('cat_pipeline',  
                                                           Pipeline(steps=[('selector',  
                                                                    Vars_selector  
                                                                    (var_idx=[0,  
                                                                    1,  
                                                                    2,  
                                                                    3,  
                                                                    4,  
                                                                    5,  
                                                                    6,  
                                                                    7,  
                                                                    8,  
                                                                    9,  
                                                                    10,  
                                                                    11,  
                                                                    12,  
                                                                    13,  
                                                                    14,  
                                                                    15,  
                                                                    16,  
                                                                    17,  
                                                                    18,  
                                                                    19,  
                                                                    20,  
                                                                    21,  
                                                                    22,  
                                                                    23,  
                                                                    24,
```



```

25,
26,
27])),
                                ('encoder',
                                Cat_vars_enco
der()))]]))],
                                ('ensemble_clf', Ensemble_clf())]]

```

Preparing a LIME explainer

Because LIME takes only float numpy array of the original dataset, to provide a human understandable explanation, it needs to know the mapping from the encoded value to original value for the categorical variables.

This is accomplished through the `categorical_names` argument. This is a dictionary of a list of values mapped to the index of the (categorical) variable column.

```
In [30]: 1 # make the categorical original values mapping
2 ## the order of the names should correspond to their (int) encoded values
3 ### This is required for LIME, because LIME uses these to index the categori
4 categorical_vals = '''Gender=Female,Male
5 Age=<12,12~40,40~65,>65
6 BMI= Underweight,Ideal,Overweight,Obese
7 Comorbidity_diabetes=F,T
8 Comorbidity_hypertension=F,T
9 Comorbidity_CAD=F,T
10 Comorbidity_pulmonary=F,T
11 Comorbidity_hepatic=F,T
12 Comorbidity_renal=F,T
13 Comorbidity_neurological=F,T
14 Comorbidity_tumor=F,T
15 Surgical_type=Abdominal,Neurosurgery,Orthopedics,Thoracic,Throat,Others
16 Emergency=F,T
17 Trauma=F,T
18 Anaesthetic_type=Local,General
19 Operative_position=Left lateral decubitus, Right lateral decubitus,Prone,Lit
20 ASA=I,II,III,IV,V
21 Timing_arrest=Induction,Intubation,Surgery,Unknown
22 Defibrillate=F,T
23 Cause_arrest=Anesthesia,Comorbidity,Surgery,Unknown
24 Hemorrhage=<200ml,200~800ml,>800ml
25 Blood_transfusion=<200ml,200~800ml,>800ml
26 Epinephrine=0,0~5mg,>5mg
27 Atropine=0,<0.65mg,>0.65mg
28 Amiodarone=0,0.04~0.3g
29 Ephedrine=0,6~15mg
30 Methoxamine=0,1~35mg
31 CPR=<30min,30~60min,>60min'''.split('\n')
32
33 categorical_names = {}
34 for val in categorical_vals:
35     categorical_names[categorical_var_idx[categorical_vars.index(val.split(
```

In [31]: 1 categorical_names

```
Out[31]: {0: ['Female', 'Male'],
1: ['<12', '12~40', '40~65', '>65 '],
2: [' Underweight', 'Ideal', 'Overweight', 'Obese'],
3: ['F', 'T'],
4: ['F', 'T'],
5: ['F', 'T'],
6: ['F', 'T'],
7: ['F', 'T'],
8: ['F', 'T'],
9: ['F', 'T'],
10: ['F', 'T'],
11: ['Abdominal',
'Neurosurgery',
'Orthopedics',
'Thoracic',
'Throat',
'Others'],
12: ['F', 'T'],
13: ['F', 'T'],
14: ['Local', 'General'],
15: ['Left lateral decubitus',
' Right lateral decubitus',
'Prone',
'Lithotomy',
'Supine'],
16: ['I', 'II', 'III', 'IV', 'V'],
17: ['Induction', 'Intubation', 'Surgery', 'Unknown'],
18: ['F', 'T'],
19: ['Anesthesia', 'Comorbidity', 'Surgery', 'Unknown'],
20: ['<200ml', '200~800ml', '>800ml '],
21: ['<200ml', '200~800ml', '>800ml '],
22: ['0', '0~5mg', '>5mg '],
23: ['0', '<0.65mg', '>0.65mg'],
24: ['0', '0.04~0.3g '],
25: ['0', '6~15mg'],
26: ['0', '1~35mg'],
27: ['<30min', '30~60min', '>60min']}
```

LIME perturbs the training sample to build a locally linear model to approximate the target model at the test point.

```
In [32]: 1 explainer = LimeTabularExplainer(
2         data_train.iloc[:, :-1].values, # remove the target variable from the tr
3         class_names=['Died', 'Survived'],
4         feature_names=list(data_train.columns),
5         categorical_features=categorical_var_idx,
6         categorical_names=categorical_names,
7         verbose=True
8     )
```

```
In [33]: 1 #https://stackoverflow.com/questions/60178732/problem-fixing-future-warning-
2
3 #test_sample = data_test.ix[data_test.Died == (not survived), :-1]
4 #test_sample = data_test.loc[data_test.Died == (not survived)][:-1]
```

```
In [34]: 1 def get_test_sample(model, died=True, correct=True, seed=42, count=1):
2     '''
3     Return one random sample from test set based on the selection criteria.
4
5     parameters:
6     model - the model which is tested for
7     survived - select survived sample
8     correct - select the sample which the model predicted correctly
9     '''
10    test_sample = data_test.ix[data_test.Died == (not died), :-1]
11    prediction = model.predict_proba(test_sample)
12
13    if (died and correct) or (not died and not correct):
14        ids = np.argwhere((prediction[:,0] > prediction[:,1]))
15    elif (died and not correct) or (not died and correct):
16        ids = np.argwhere((prediction[:,0] <= prediction[:,1]))
17    if count == 1:
18        idx = random.Random(seed).choice(ids)
19    else:
20        ids = ids.reshape(len(ids)).tolist()
21        idx = random.Random(seed).sample(ids, count)
22
23    return test_sample.iloc[idx,:]
```

```
In [35]: 1 def plot_vertical_bar(explanation, title='Local explanation for class died'
2     exp_list = explanation.as_list()
3     tags, values, colors = [],[],[]
4     for i in range(len(exp_list)):
5         tags.append(exp_list[i][0])
6         values.append(exp_list[i][1])
7         if(exp_list[i][1] > 0):
8             colors.append('r')
9         else:
10            colors.append('g')
11
12    fig, ax = plt.subplots()
13
14    y_pos = np.arange(len(tags))
15
16    ax.barh(y_pos, values, align='center',
17            color=colors, ecolor='black')
18    ax.set_yticks(y_pos)
19    ax.set_yticklabels(tags)
20    ax.invert_yaxis() # Labels read top-to-bottom
21    ax.set_title(title)
```

```
In [36]: 1 def score_compare(x, y, threshold = 0.5):
2         if x < threshold and y < threshold:
3             return 'True Negative'
4         elif x < threshold and y > threshold:
5             return 'False Negative'
6         elif x > threshold and y > threshold:
7             return 'True Positive'
8         elif x > threshold and y < threshold:
9             return 'False Positive'
```

Correct, True Negative, patient survived

```
In [37]: 1 test_sample_died_correct = get_test_sample(ensemble_pipeline, died=True, cor
2         test_sample_died_correct
```

.ix is deprecated. Please use
.loc for label based indexing or
.iloc for positional indexing

See the documentation here:

http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated (http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated)

.ix is deprecated. Please use
.loc for label based indexing or
.iloc for positional indexing

See the documentation here:

http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated (http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated)

Out[37]:

	Gender	Age	BMI	Comorbidity_diabetes	Comorbidity_hypertension	Comorbidity_CAD	Comorbidity
17	0	3	2	0	1	1	

```
In [38]: 1 # prediction
2         result = ensemble_pipeline.predict_proba(test_sample_died_correct)
3         score = 1-round(result[0,0],2)
4         score
```

Out[38]: 0.19999999999999996

Let's see how LIME explainer (with default settings) explains it.

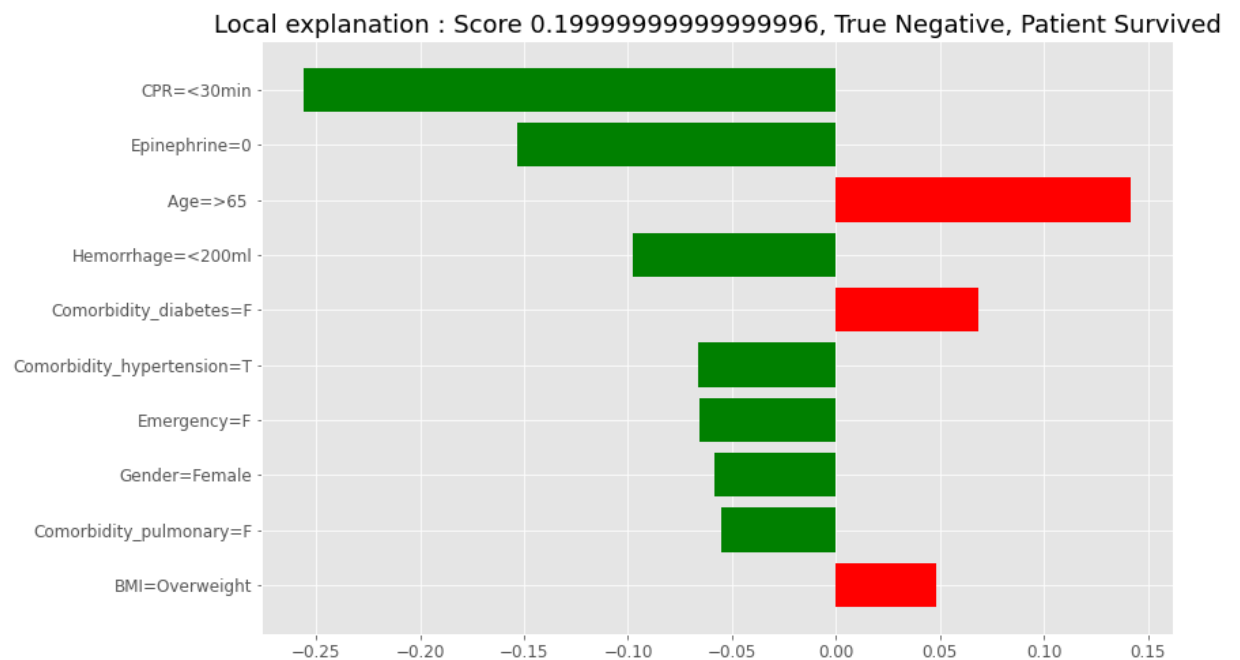
```
In [39]: 1 explanation = explainer.explain_instance(test_sample_died_correct.values.ra
```

Intercept 0.7589152534734966

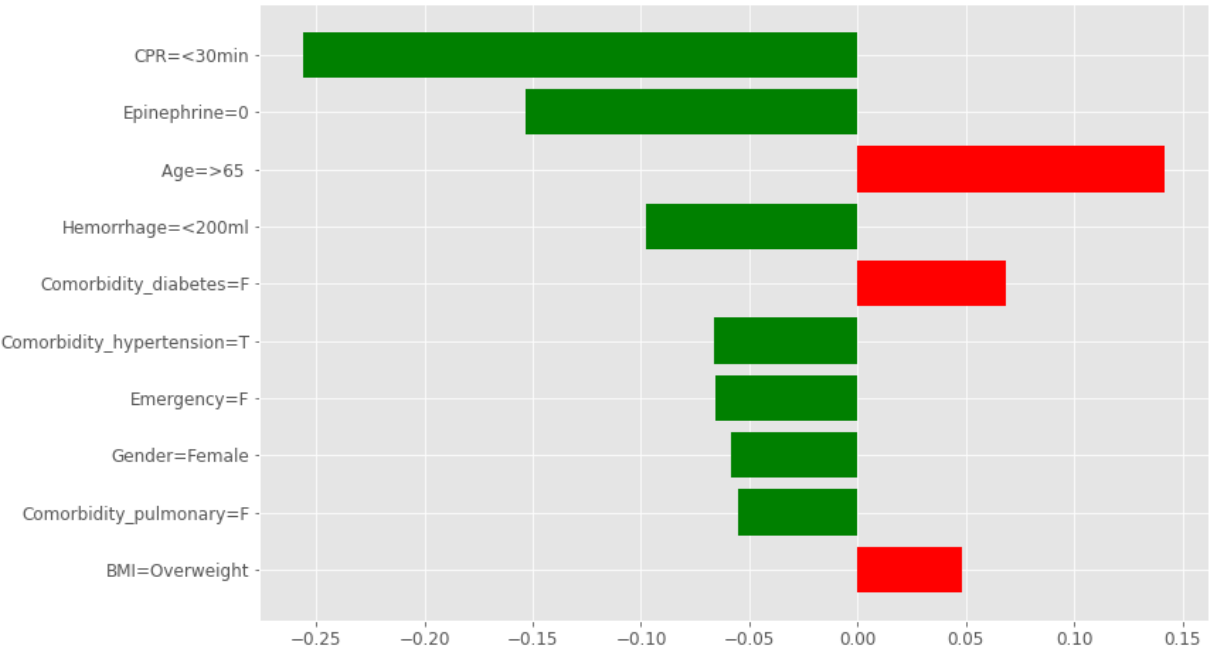
Prediction_local [0.26636104]

Right: 0.20438393533185692

```
In [40]: 1 outcome = 0
2 outcome_title = 'Survived'
3 plot_vertical_bar(explanation, 'Local explanation : Score {}'.format(outcome), Patient {
```



```
In [41]: 1 outcome = 0
2 outcome_title = 'Survived'
3 #plot_vertical_bar(explanation, 'Local explanation : Score {}'.format(score), {}, Patient
4 plot_vertical_bar(explanation, '{}.format(score, score_compare(score, outcome_title))
```



```
In [42]: 1 test_sample_died_correct = get_test_sample(ensemble_pipeline, seed=43) #, su
2 test_sample_died_correct
```

.ix is deprecated. Please use
.loc for label based indexing or
.iloc for positional indexing

See the documentation here:
http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated (http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated)

.ix is deprecated. Please use
.loc for label based indexing or
.iloc for positional indexing

See the documentation here:
http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated (http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated)

Out[42]:

	Gender	Age	BMI	Comorbidity_diabetes	Comorbidity_hypertension	Comorbidity_CAD	Comor
139	0	2	0	0	0	1	

```
In [43]: 1 # prediction
2 result = ensemble_pipeline.predict_proba(test_sample_died_correct)
3 score = 1-round(result[0,0],2)
4 score
```

Out[43]: 0.26

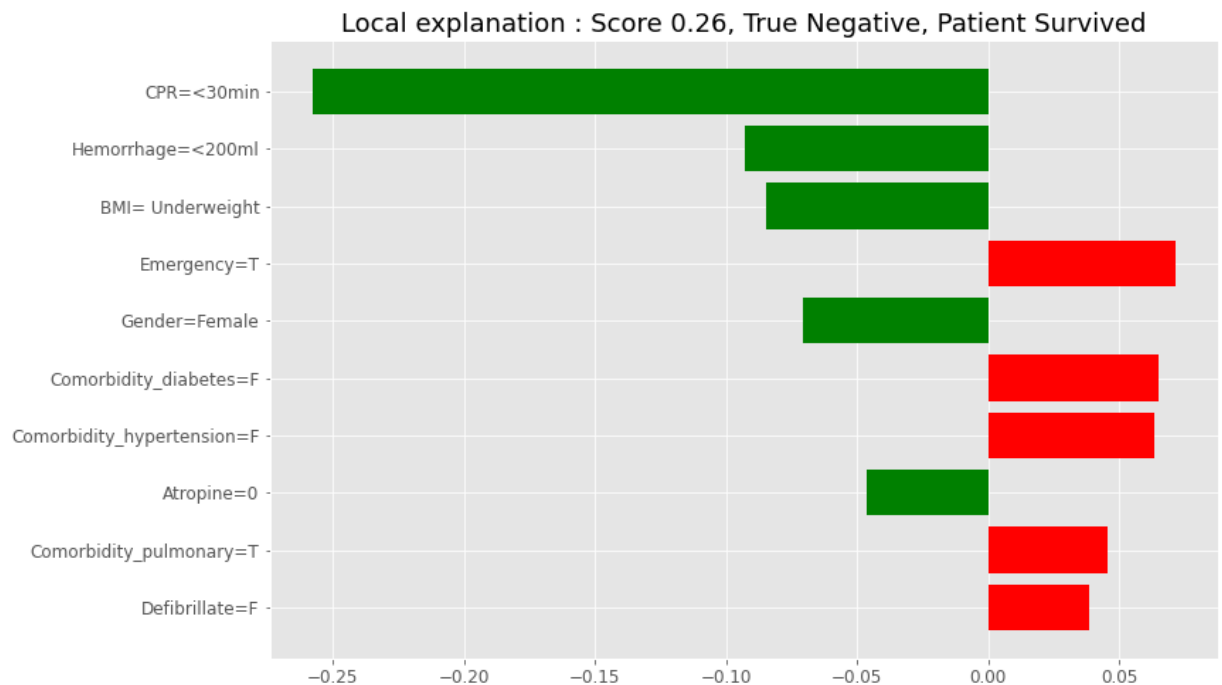
```
In [44]: 1 explanation = explainer.explain_instance(test_sample_died_correct.values.ra
2 print(explanation)
3 outcome = 0
4 outcome_title = 'Survived'
5 plot_vertical_bar(explanation, 'Local explanation : Score {}, {}, Patient {'
```

Intercept 0.6035533371112072

Prediction_local [0.33570023]

Right: 0.26366191167685205

<lime.explanation.Explanation object at 0x000000C8B3A919E8>



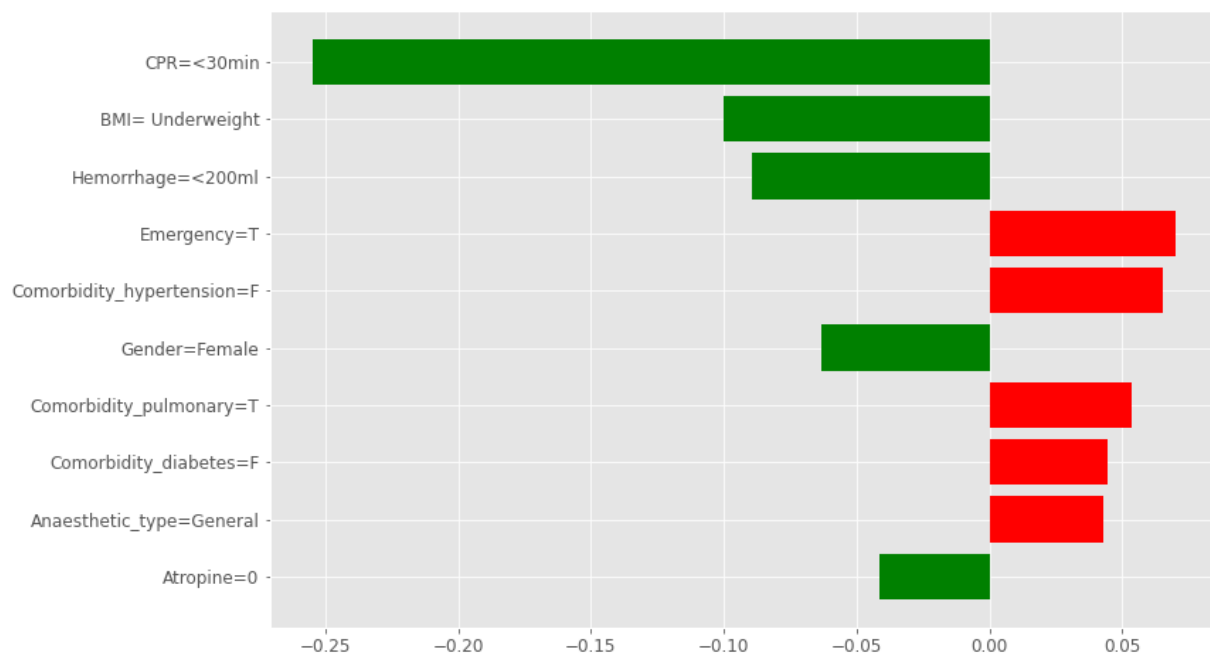

```
In [45]: 1 explanation = explainer.explain_instance(test_sample_died_correct.values.ra
2 print(explanation)
3 outcome = 0
4 outcome_title = 'Survived'
5 plot_vertical_bar(explanation, '{}.format(score, score_compare(score, outcom
```

Intercept 0.5931771988539475

Prediction_local [0.32067057]

Right: 0.26366191167685205

<lime.explanation.Explanation object at 0x000000C8B4EFF9B0>



Correct, True Positive, patient died

```
In [46]: 1 test_sample_died_correct = get_test_sample(ensemble_pipeline, died=False, co
2 test_sample_died_correct
```

.ix is deprecated. Please use
 .loc for label based indexing or
 .iloc for positional indexing

See the documentation here:

http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated (http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated)

.ix is deprecated. Please use
 .loc for label based indexing or
 .iloc for positional indexing

See the documentation here:

http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated (http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated)

Out[46]:

	Gender	Age	BMI	Comorbidity_diabetes	Comorbidity_hypertension	Comorbidity_CAD	Comor
104	0	1	3	0	0	1	

```
In [47]: 1 # prediction
2 result = ensemble_pipeline.predict_proba(test_sample_died_correct)
3 score = 1-round(result[0,0],2)
4 score
```

Out[47]: 0.9

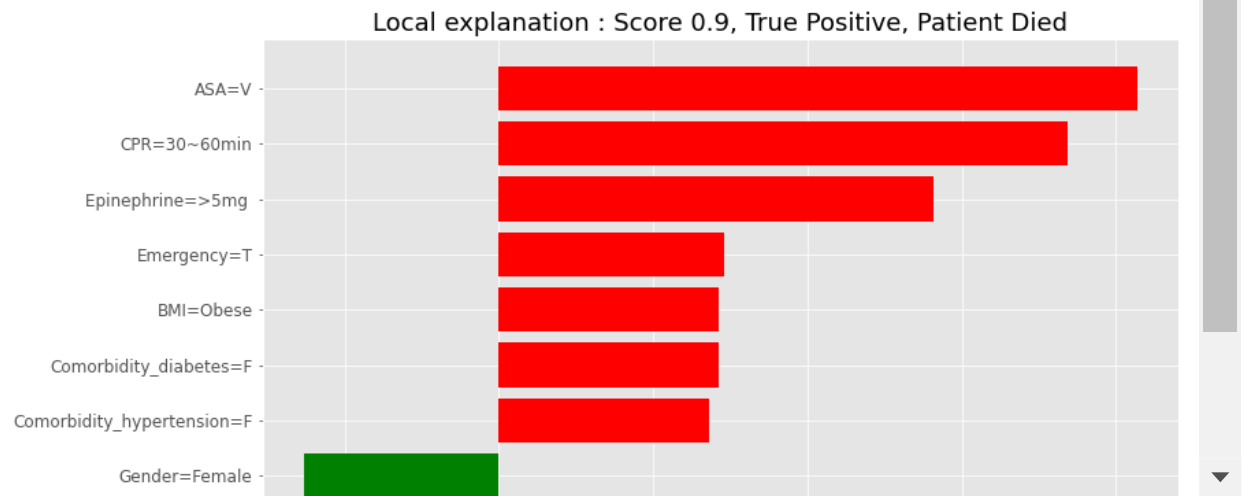
```
In [48]: 1 explanation = explainer.explain_instance(test_sample_died_correct.values.ra
2 print(explanation)
3 outcome = 1
4 outcome_title = 'Died'
5 plot_vertical_bar(explanation, 'Local explanation : Score {}, {}, Patient {'
```

Intercept 0.25696037453462417

Prediction_local [1.00166513]

Right: 0.8997209131124573

<lime.explanation.Explanation object at 0x000000C8B24E2160>



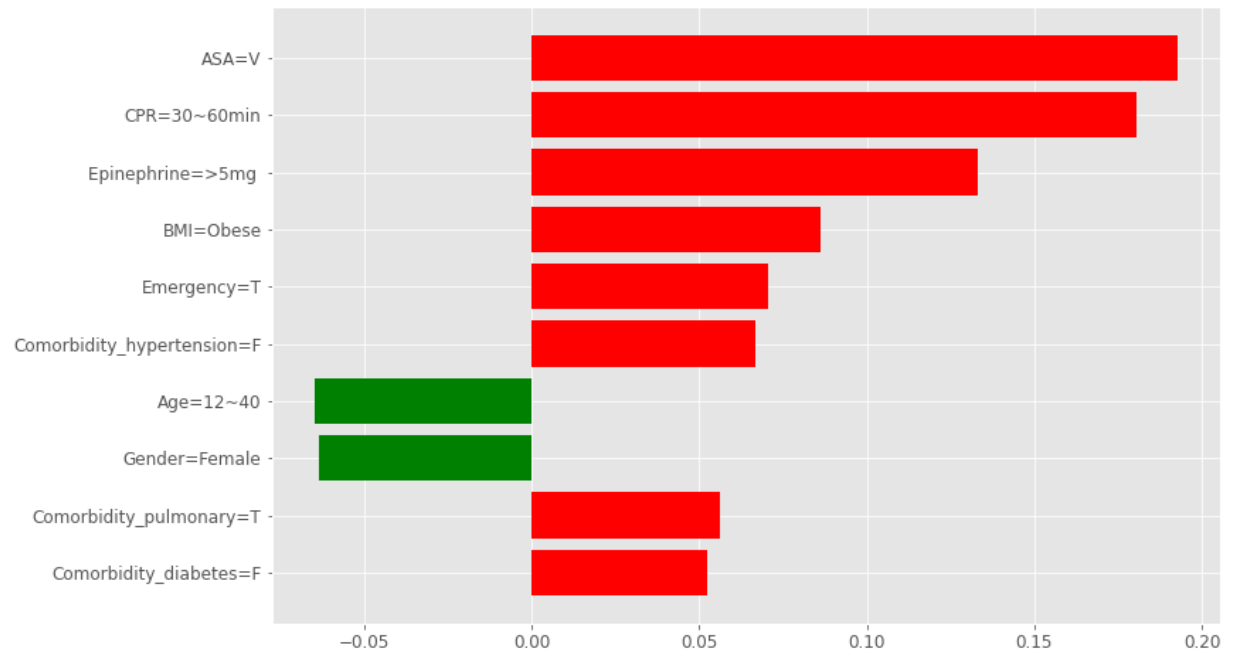
```
In [49]: 1 explanation = explainer.explain_instance(test_sample_died_correct.values.ra
2 print(explanation)
3 outcome = 1
4 outcome_title = 'Died'
5 plot_vertical_bar(explanation, '{}.format(score, score_compare(score, outcom
```

Intercept 0.2808745369165182

Prediction_local [0.99104523]

Right: 0.8997209131124573

<lime.explanation.Explanation object at 0x000000C8B4DE6978>



```
In [50]: 1 test_sample_died_correct = get_test_sample(ensemble_pipeline, died=False, se
2 test_sample_died_correct
```

.ix is deprecated. Please use
.loc for label based indexing or
.iloc for positional indexing

See the documentation here:

http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated (http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated)

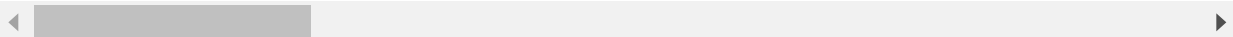
.ix is deprecated. Please use
.loc for label based indexing or
.iloc for positional indexing

See the documentation here:

http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated (http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated)

Out[50]:

	Gender	Age	BMI	Comorbidity_diabetes	Comorbidity_hypertension	Comorbidity_CAD	Comor
100	0	1	2	0	1	1	



```
In [51]: 1 # prediction
2 result = ensemble_pipeline.predict_proba(test_sample_died_correct)
3 score = 1-round(result[0,0],2)
4 score
```

Out[51]: 0.89

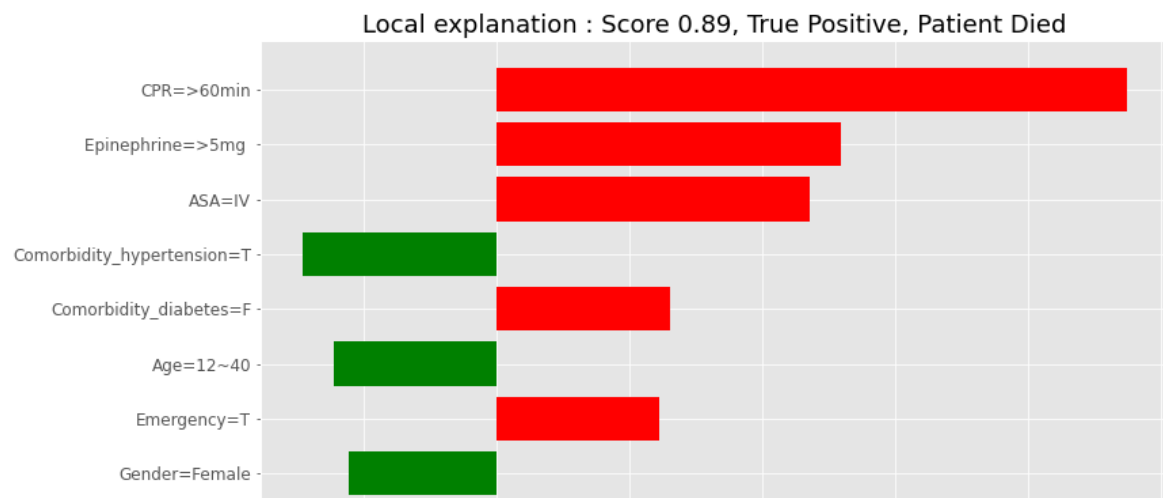
```
In [52]: 1 explanation = explainer.explain_instance(test_sample_died_correct.values.ra
2 print(explanation)
3 outcome = 1
4 outcome_title = 'Died'
5 plot_vertical_bar(explanation, 'Local explanation : Score {}, {}, Patient {'
```

Intercept 0.3926734633053218

Prediction_local [0.90741066]

Right: 0.8863784910412932

<lime.explanation.Explanation object at 0x000000C8B4E7D588>



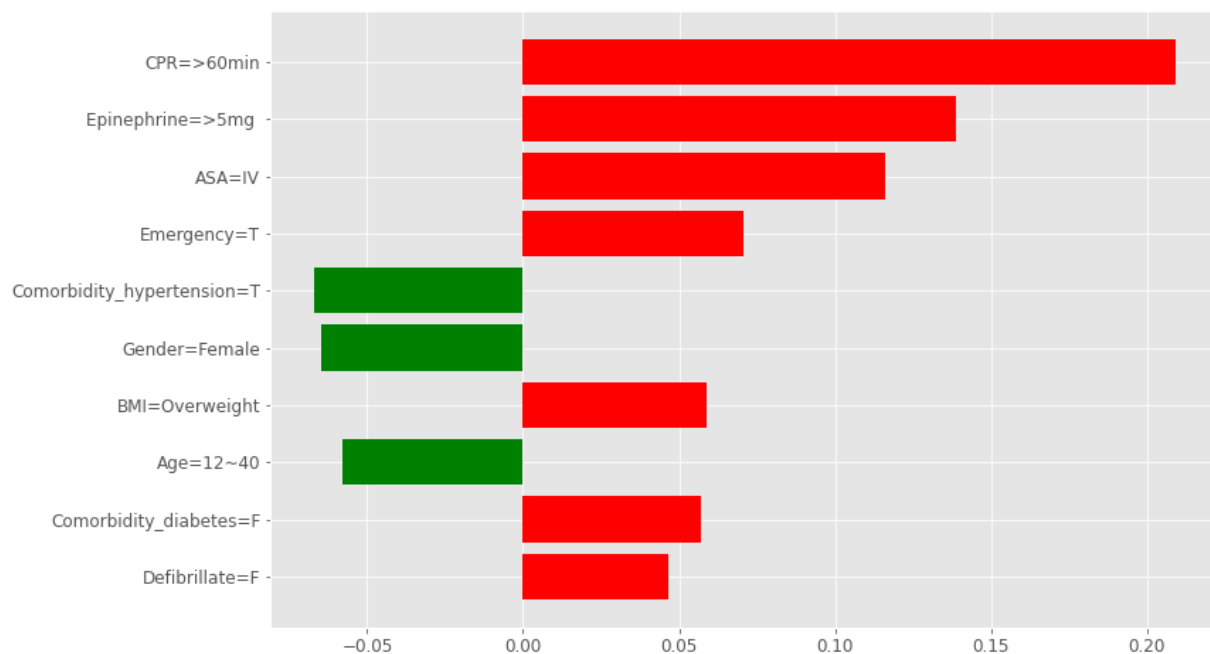
```
In [53]: 1 explanation = explainer.explain_instance(test_sample_died_correct.values.ra
2 print(explanation)
3 outcome = 0
4 outcome_title = 'survived'
5 plot_vertical_bar(explanation, ''.format(score, score_compare(score, outcom
```

Intercept 0.3841302201029104

Prediction_local [0.89185997]

Right: 0.8863784910412932

<lime.explanation.Explanation object at 0x000000C8B4E3DA58>



Incorrect, False Negative, patient died

```
In [54]: 1 test_sample_died_correct = get_test_sample(ensemble_pipeline, died=False, co
2 test_sample_incorrect
```

.ix is deprecated. Please use
.loc for label based indexing or
.iloc for positional indexing

See the documentation here:

http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated (http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated)

.ix is deprecated. Please use
.loc for label based indexing or
.iloc for positional indexing

See the documentation here:

http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated (http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated)

```
-----
NameError                                Traceback (most recent call last)
<ipython-input-54-bdb0d0cc47fa> in <module>
      1 test_sample_died_correct = get_test_sample(ensemble_pipeline, died=False, correct=False, seed=42, count=1)
----> 2 test_sample_incorrect

NameError: name 'test_sample_incorrect' is not defined
```

```
In [ ]: 1 # prediction
2 result = ensemble_pipeline.predict_proba(test_sample_incorrect)
3 score = 1-round(result[0,0],2)
4 score
```

```
In [ ]: 1 explanation = explainer.explain_instance(test_sample_incorrect.values.ravel()
2 print(explanation)
3 outcome = 1
4 outcome_title = 'Died'
5 plot_vertical_bar(explanation, 'Local explanation : Score {}, {}, Patient {
```

```
In [ ]: 1 explanation = explainer.explain_instance(test_sample_incorrect.values.ravel()
2 print(explanation)
3 outcome = 1
4 outcome_title = 'Died'
5 plot_vertical_bar(explanation, ''.format(score, score_compare(score, outcom
```

```
In [ ]: 1 test_sample_died_correct = get_test_sample(ensemble_pipeline, died=False, co
2 test_sample_incorrect
```



```
In [ ]: 1 # prediction
        2 result = ensemble_pipeline.predict_proba(test_sample_incorrect)
        3 score = 1-round(result[0,0],2)
        4 score
```

```
In [ ]: 1 explanation = explainer.explain_instance(test_sample_incorrect.values.ravel
        2 print(explanation)
        3 outcome = 1
        4 outcome_title = 'Died'
        5 plot_vertical_bar(explanation, 'Local explanation : Score {}, {}, Patient {
```

Incorrect, False Positive, patient survived

```
In [55]: 1 test_sample_incorrect = get_test_sample(ensemble_pipeline, correct=False, di
        2 test_sample_incorrect
```

.ix is deprecated. Please use
.loc for label based indexing or
.iloc for positional indexing

See the documentation here:

http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated (http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated)

.ix is deprecated. Please use
.loc for label based indexing or
.iloc for positional indexing

See the documentation here:

http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated (http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated)

Out[55]:

	Gender	Age	BMI	Comorbidity_diabetes	Comorbidity_hypertension	Comorbidity_CAD	Comorbidity_CAD
19	1	2	3	0	0	1	

```
In [56]: 1 # prediction
        2 result = ensemble_pipeline.predict_proba(test_sample_incorrect)
        3 score = 1-round(result[0,0],2)
        4 score
```

Out[56]: 0.71

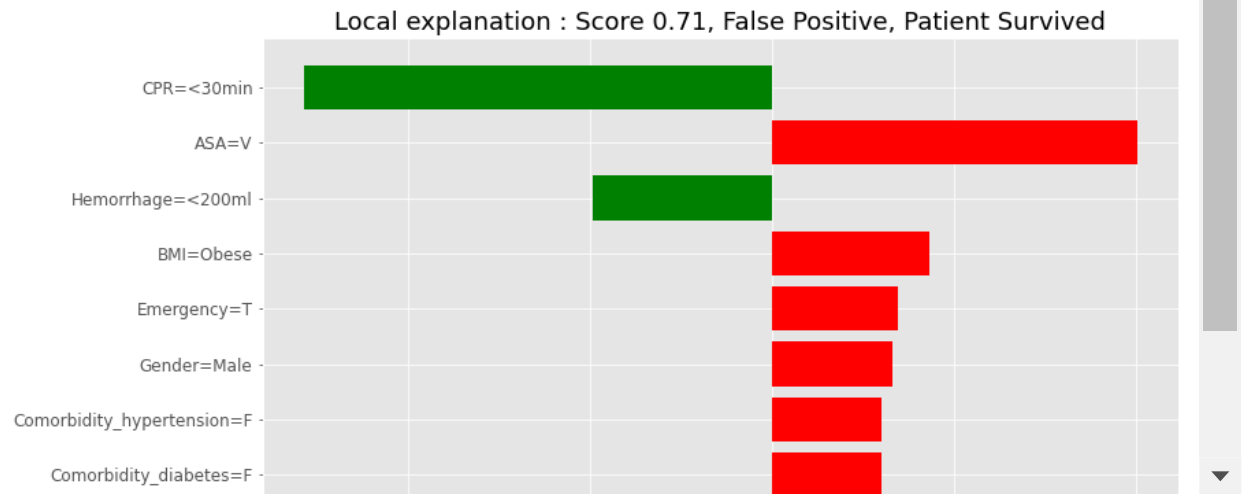
```
In [57]: 1 explanation = explainer.explain_instance(test_sample_incorrect.values.ravel  
2 print(explanation)  
3 outcome = 0  
4 outcome_title = 'Survived'  
5 plot_vertical_bar(explanation, 'Local explanation : Score {}, {}, Patient {
```

Intercept 0.5335559116085485

Prediction_local [0.72396846]

Right: 0.7144154525359052

<lime.explanation.Explanation object at 0x000000C8B39BBD68>



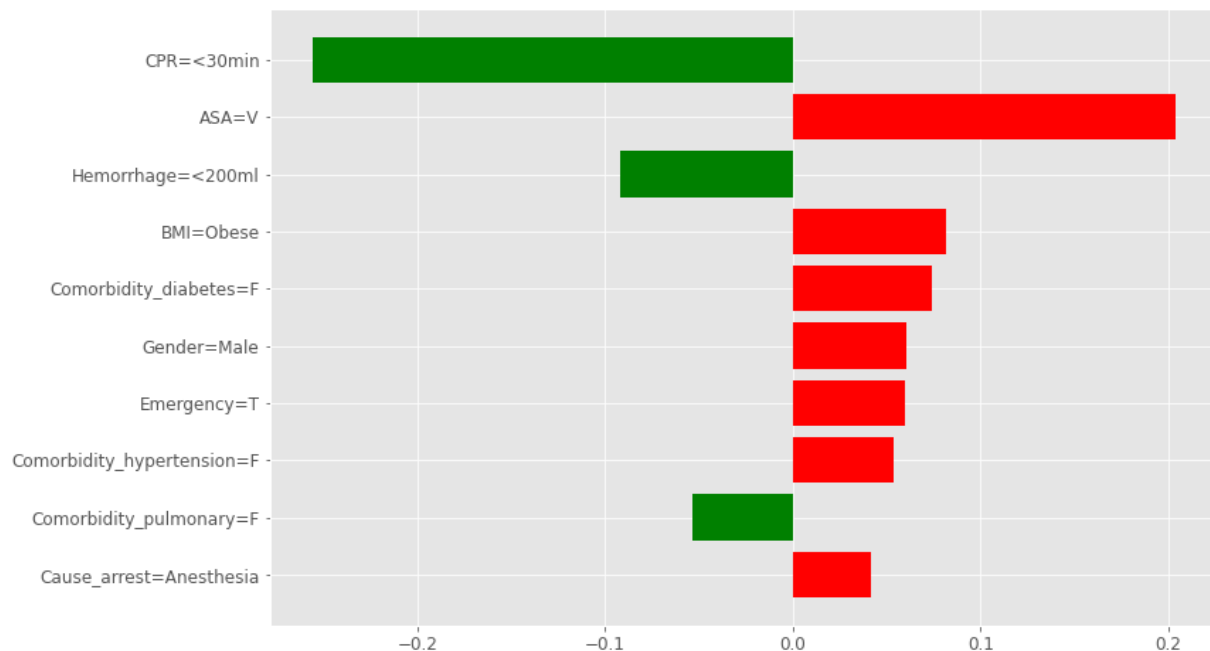
```
In [58]: 1 explanation = explainer.explain_instance(test_sample_incorrect.values.ravel
2 print(explanation)
3 outcome = 0
4 outcome_title = 'Survived'
5 #plot_vertical_bar(explanation, 'Local explanation : Score {}, {}, Patient
6 plot_vertical_bar(explanation, '{}.format(score, score_compare(score, outcom
```

Intercept 0.5379357467805043

Prediction_local [0.71395739]

Right: 0.7144154525359052

<lime.explanation.Explanation object at 0x000000C8B391D978>



```
In [59]: 1 test_sample_incorrect = get_test_sample(ensemble_pipeline, correct=False, di
2 test_sample_incorrect
```

.ix is deprecated. Please use
.loc for label based indexing or
.iloc for positional indexing

See the documentation here:

http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated (http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated)

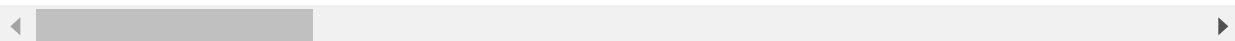
.ix is deprecated. Please use
.loc for label based indexing or
.iloc for positional indexing

See the documentation here:

http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated (http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#ix-indexer-is-deprecated)

Out[59]:

	Gender	Age	BMI	Comorbidity_diabetes	Comorbidity_hypertension	Comorbidity_CAD	Comor
126	1	2	3	0	1	1	



```
In [60]: 1 # prediction
2 result = ensemble_pipeline.predict_proba(test_sample_incorrect)
3 score = 1-round(result[0,0],2)
4 score
```

Out[60]: 0.62

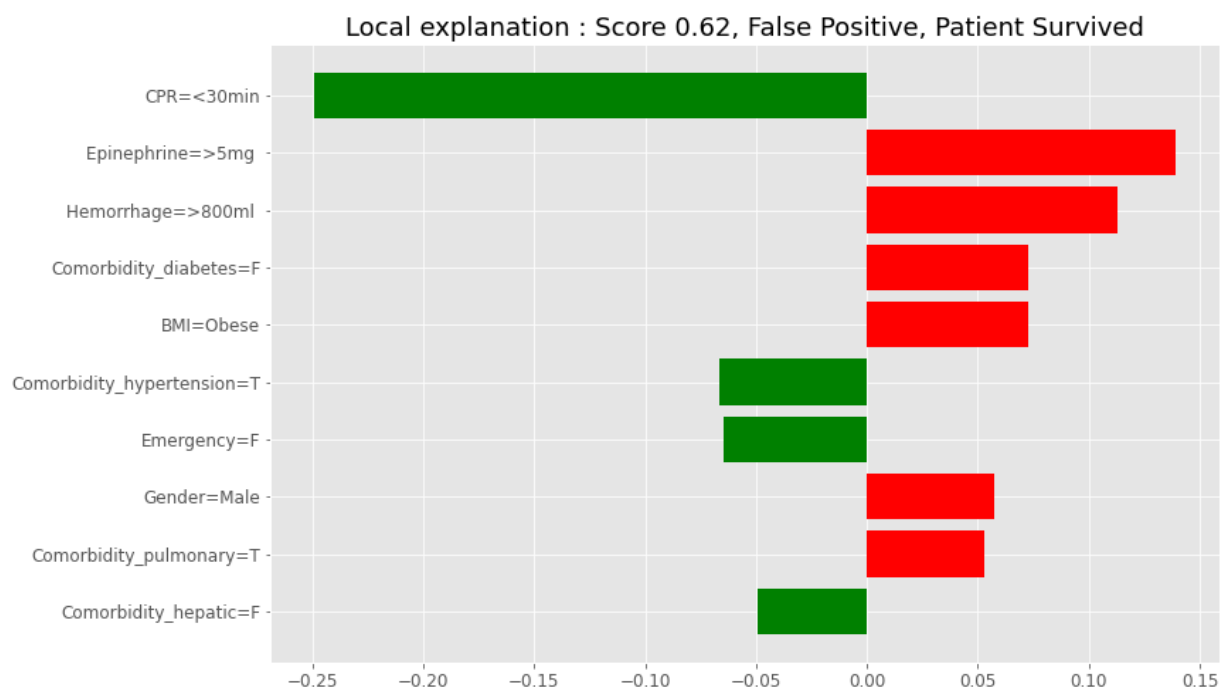
```
In [61]: 1 explanation = explainer.explain_instance(test_sample_incorrect.values.ravel
2 print(explanation)
3 outcome = 0
4 outcome_title = 'Survived'
5 plot_vertical_bar(explanation, 'Local explanation : Score {}, {}, Patient {'
```

Intercept 0.5646545897255475

Prediction_local [0.64332151]

Right: 0.6166866999359815

<lime.explanation.Explanation object at 0x000000C8B24B3518>



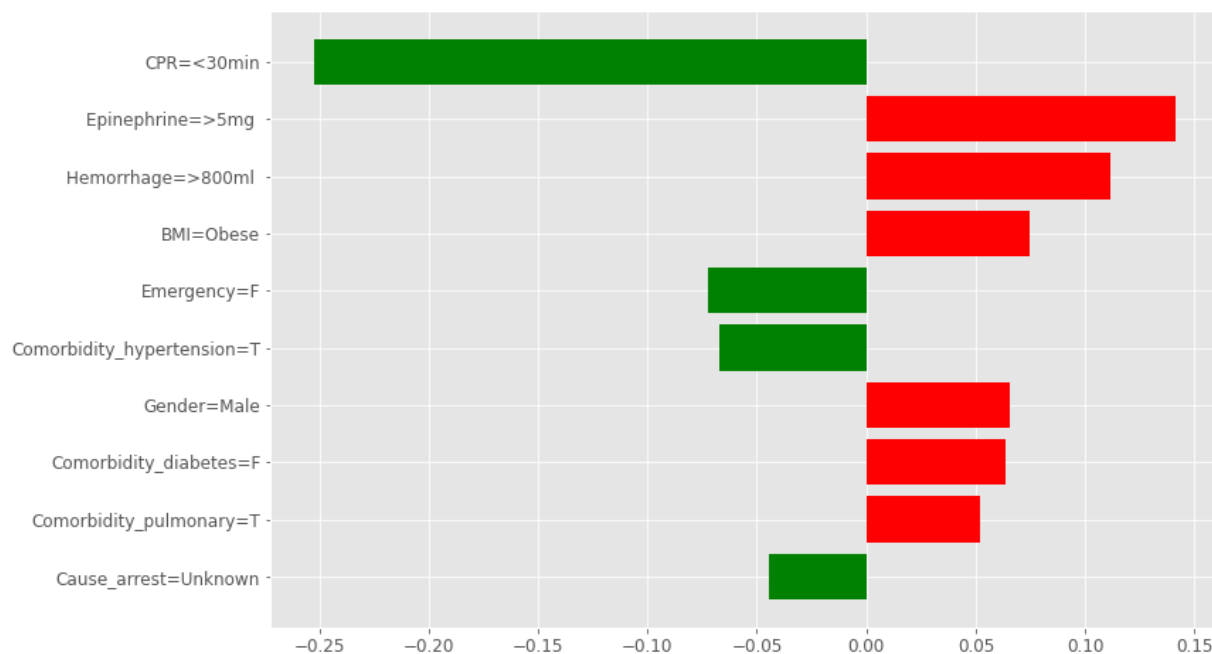
```
In [62]: 1 explanation = explainer.explain_instance(test_sample_incorrect.values.ravel
2         print(explanation)
3         outcome = 1
4         outcome_title = 'died'
5         plot_vertical_bar(explanation, '{}.format(score, score_compare(score, outcom
```

Intercept 0.5442164456791525

Prediction_local [0.61693366]

Right: 0.6166866999359815

<lime.explanation.Explanation object at 0x000000C8B20F66D8>



Overall explanation by LIME

```
In [63]: 1 #model =rf_clf
2 model=ensemble_clf
3
4 RANDOM_STATE=123
5
6 explainer = LimeTabularExplainer(
7     data_train.iloc[:, :-1].values, # remove the target variable from the tr
8     class_names=['survived','died'],
9     feature_names=list(data_train.columns),
10    categorical_features=categorical_var_idx,
11    categorical_names=categorical_names,
12    mode="classification",
13    verbose=True,
14    random_state=RANDOM_STATE)
15
16
17 ##LIME
18
19 #Explain samples in test set
20 X_explain = data_test.iloc[:, :-1]
21 #data_test.drop('Died',1)
22
23 #Explaining first subject in test set using all 30 features
24 exp = explainer.explain_instance(X_explain.values[0,:],model.predict_proba,
25 #Plot local explanation
26 plt = exp.as_pyplot_figure()
27 plt.tight_layout()
28 exp.show_in_notebook(show_table=True)
```

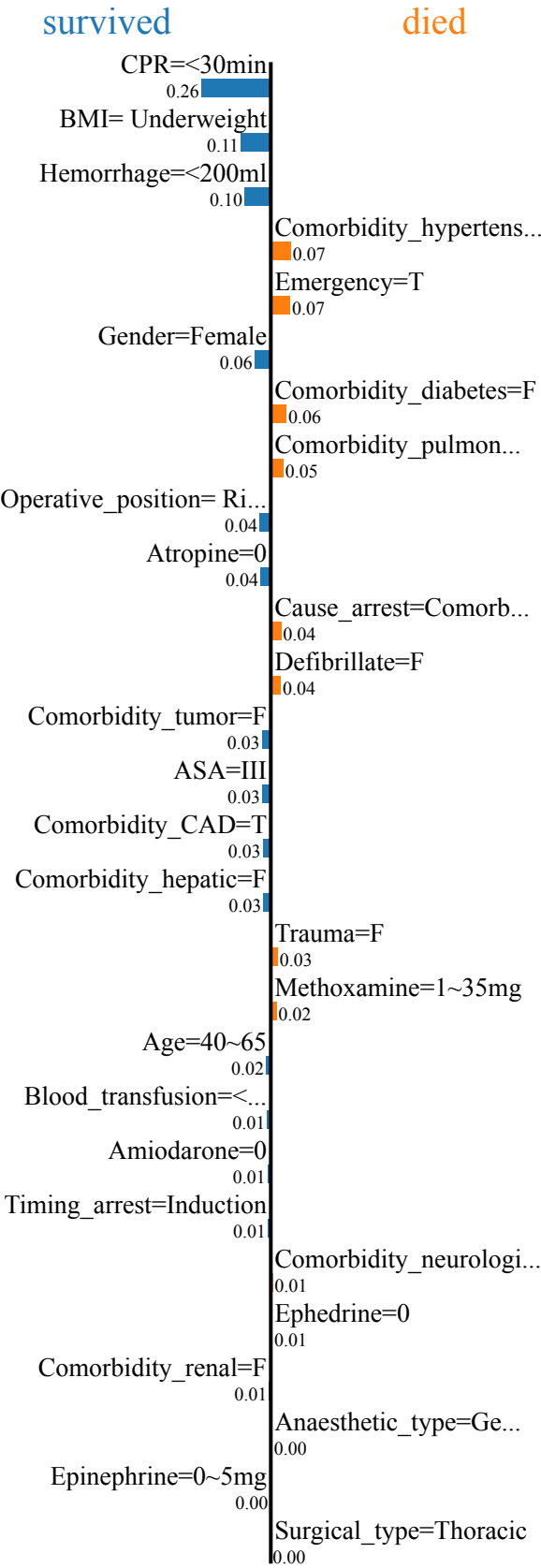
Intercept 0.6523716191933714

Prediction_local [0.25479006]

Right: 0.2636619116768521

Prediction probabilities

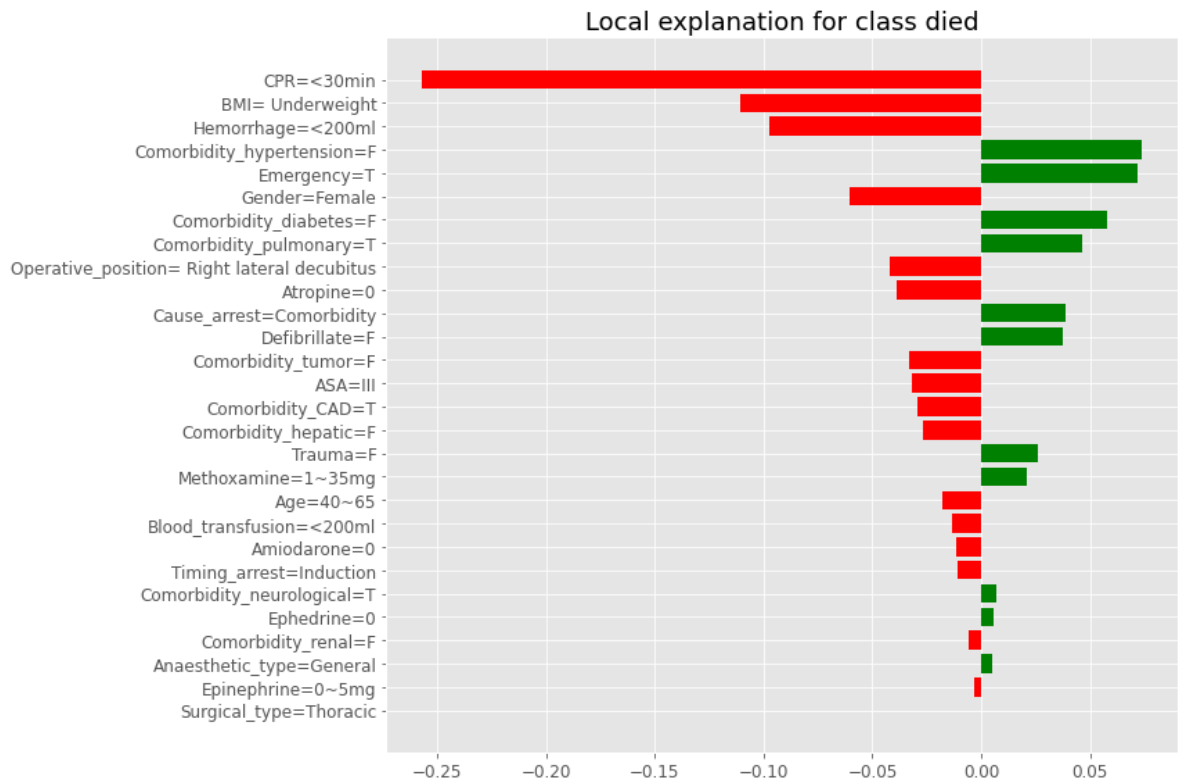




Feature Value

CPR=<30min	True
BMI= Underweight	True
Hemorrhage=<200ml	True

Comorbidity_hypertension=F	True
Emergency=T	True
Gender=Female	True
Comorbidity_diabetes=F	True
Comorbidity_pulmonary=T	True



3. Preparing a SHAP explainer

```
In [64]: 1 with open(r'C:\Users\neo\Huijie\Code\Predicting CA Mortality\Models\rf_clf_f
2         rf_clf = pickle.load(f)
3         model=rf_clf
```

Trying to unpickle estimator DecisionTreeClassifier from version 0.24.0 when using version 1.0. This might lead to breaking code or invalid results. Use at your own risk. For more info please refer to:

https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations (https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations)

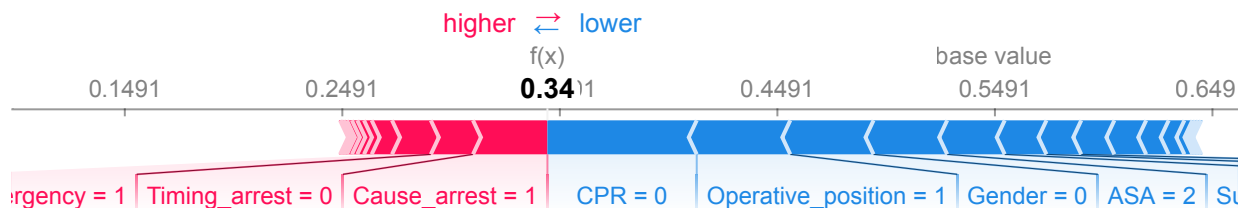
Trying to unpickle estimator RandomForestClassifier from version 0.24.0 when using version 1.0. This might lead to breaking code or invalid results. Use at your own risk. For more info please refer to:

https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations (https://scikit-learn.org/stable/modules/model_persistence.html#security-maintainability-limitations)

```
In [65]: 1 #We can also use Shapley force plots to give us additional information on a
2
3 # explain the model's predictions on test set using SHAP values
4 # same syntax works for xgboost, LightGBM, CatBoost, and some scikit-Learn m
5 explainer = shap.TreeExplainer(model)
6
7 # shap_values consists of a list of two matrices of dimension samplesize x #
8 # The first matrix uses average nr of benign samples as base value
9 # The second matrix which is used below uses average nr of malignant samples
10 shap_values = explainer.shap_values(X_explain)
11
12
13 # Interactive visualization of the explanation of the first subject
14 # in the test set (X_explain).
15 # It shows the relative contribution of features to get from the base value
16 # samples(average value)
17 # to the output value (1 in case of malignant sample)
18 # the numbers at the bottom show the actual values for this sample.
19 shap.initjs() #initialize javascript in cell
20 shap.force_plot(explainer.expected_value[1], shap_values[1][0,:], X_explain.
```

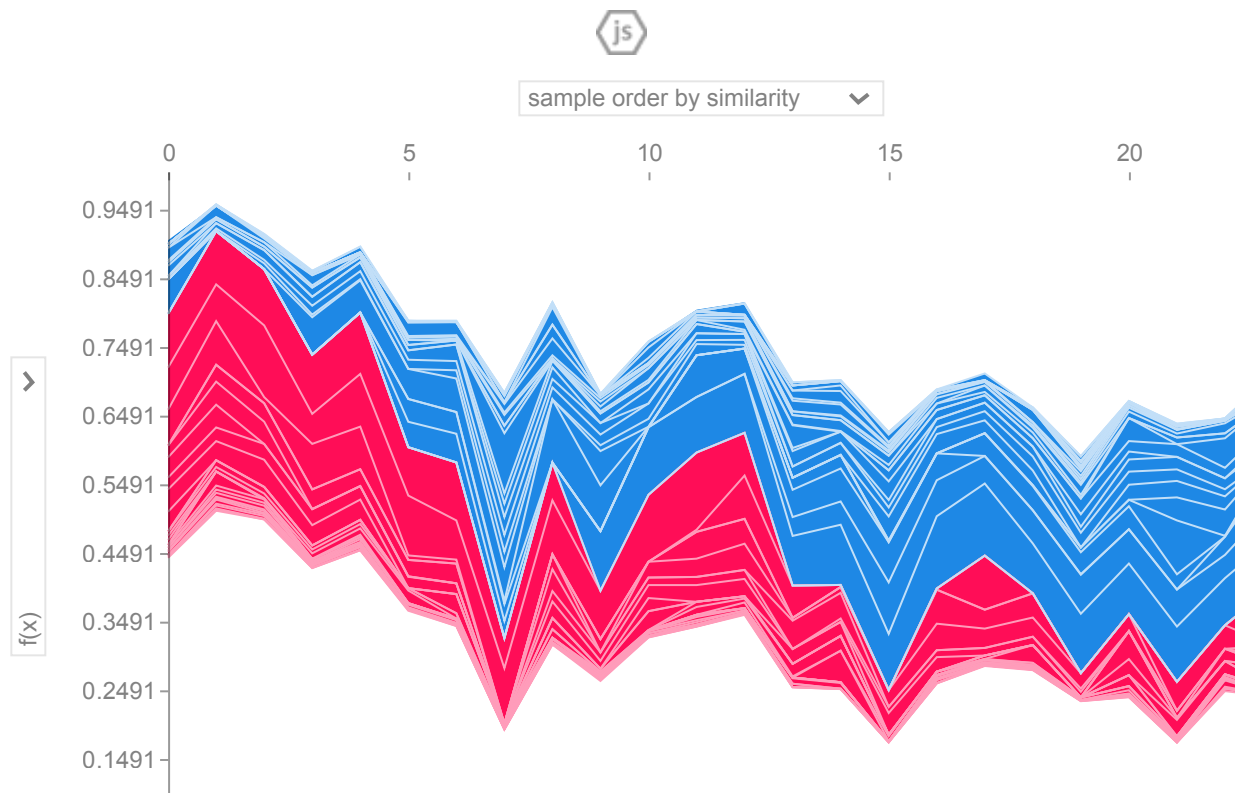


Out[65]:



```
In [66]: 1 #Interactive visualization of all sample/feature Shapley values  
2 #It is possible to show the relative contribution of individual features for  
3 # samples on the y-axis as well.  
4 shap.initjs()  
5 shap.force_plot(explainer.expected_value[1], shap_values[1], X_explain)
```

Out[66]:



In [67]:

1

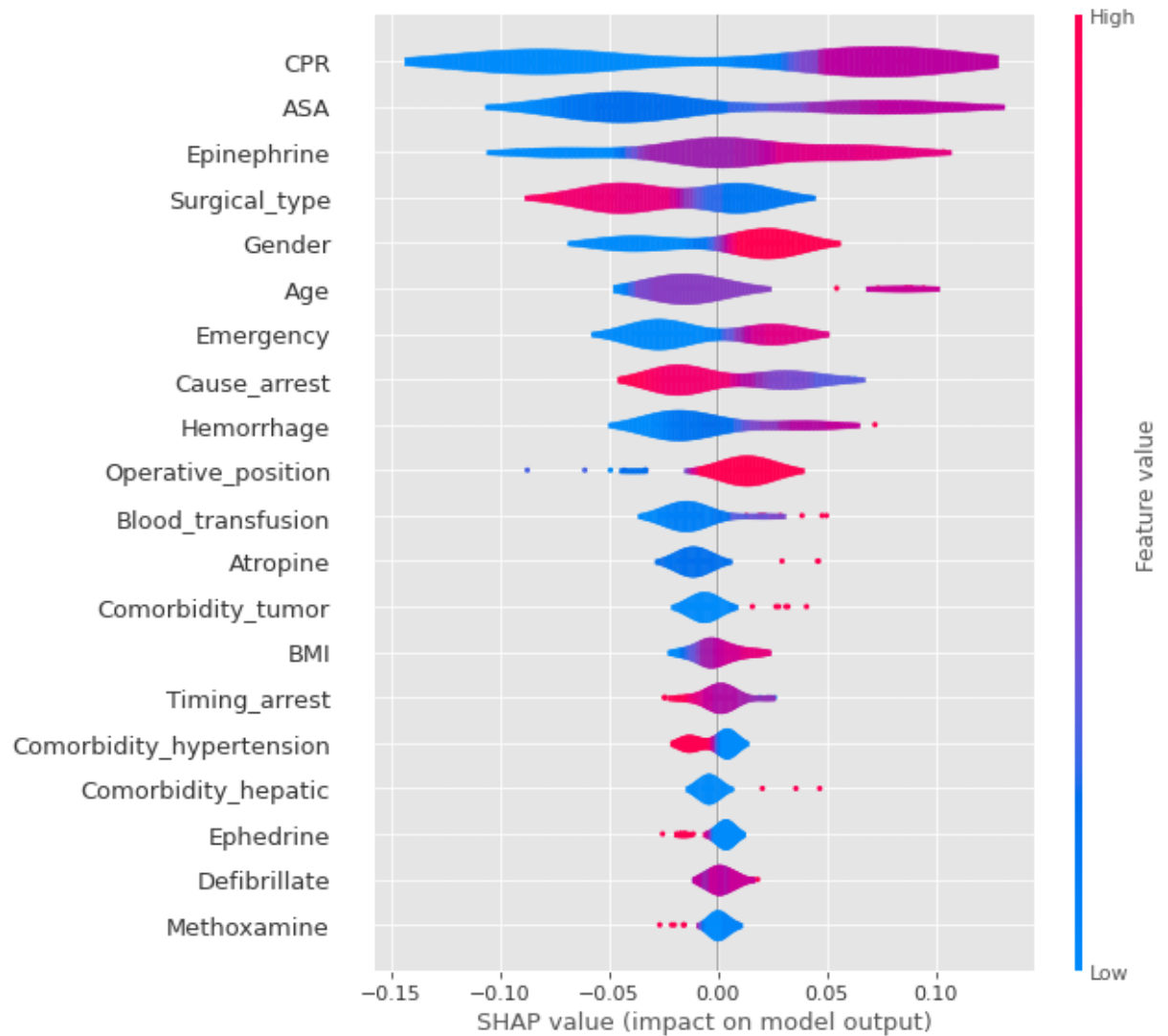
```
#A summary plot with the shapley value (feature importance)
```

2

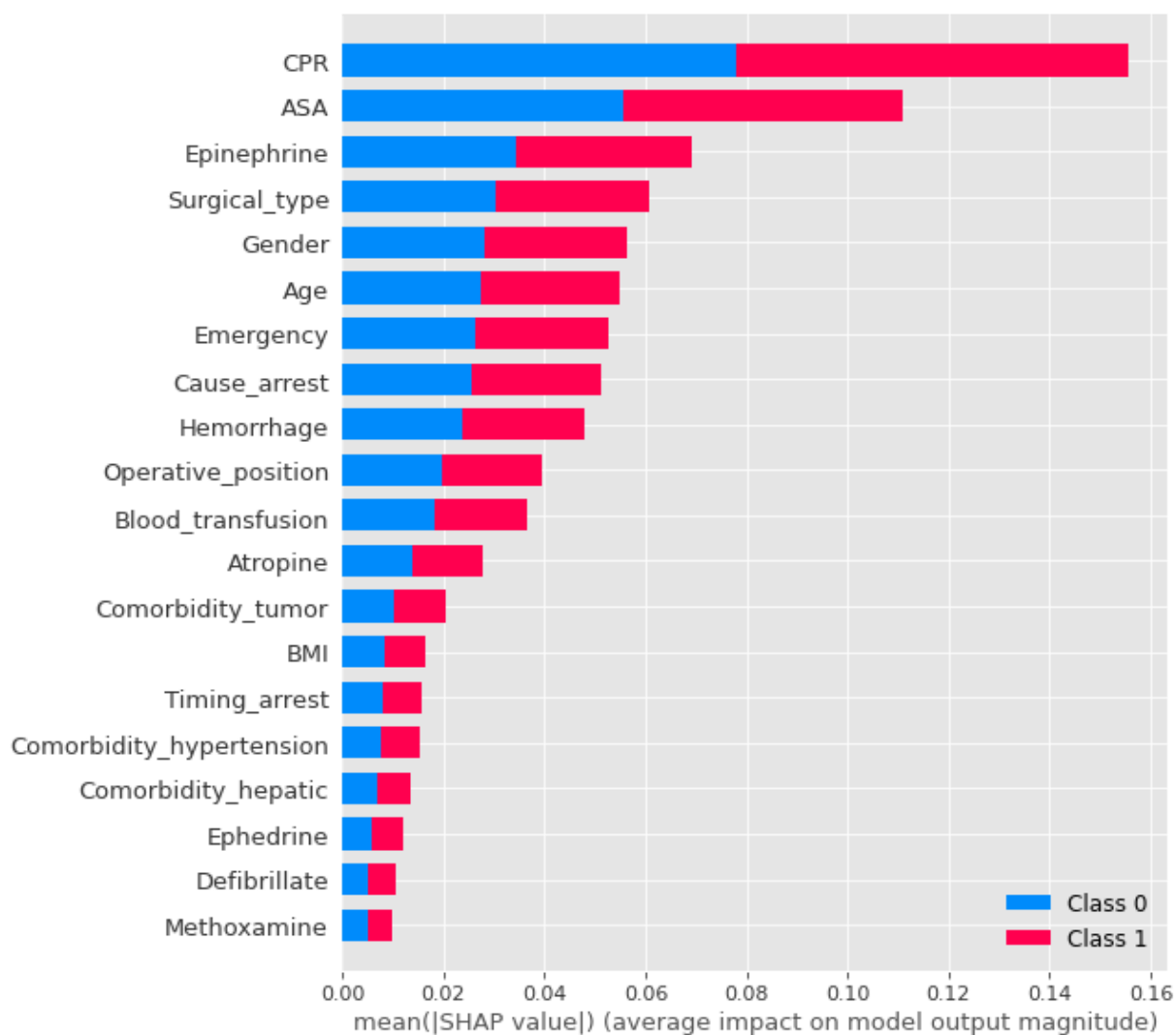
```
shap.summary_plot(shap_values[1], X_explain)
```



```
In [68]: 1 #Same as above, but with violin plots to better see the distribution of shap
2 shap.summary_plot(shap_values[1], X_explain, plot_type="violin")
```



```
In [69]: 1 shap.summary_plot(shap_values, X_explain, plot_type="bar", show=False)
2 plt.tight_layout()
3 plt.savefig('test.png', dpi=300)
```



4. Partial Dependence Plots

Partial Dependency Plots (DPD) show the effect a feature has on the outcome of a predictive based model. It marginalizes the model output over the distribution of features in order to extract the importance of the feature of interest. This importance calculation is based on an important assumption, namely that the feature of interest is not correlated with all other features (except for the target). The reason for this is that it will show data points that are likely to be impossible. For example, weight and height are correlated but the PDP might show the effect of a large weight and very small height on the target while that combination is highly unlikely. This can be partially resolved by showing a rug at the bottom of your PDP.

<https://lmc2179.github.io/posts/pdp.html> (<https://lmc2179.github.io/posts/pdp.html>)

```
In [70]: 1 from pdpbox import pdp
        2 import seaborn as sns
```

```
In [71]: 1 #data=data_train.drop('Died',1)
        2 df_pdp=data_train_X
```

```
In [72]: 1 with open(r'C:\Users\neo\Huijie\Code\Predicting CA Mortality\Models\rf_clf_f
        2         rf_clf = pickle.load(f)
        3 model=rf_clf
```

4.0 Correlation - One-hot Encoding

Below the correlation matrix is shown between features to give an indication of whether the assumption of independence is violated or not. From these features one can conclude that there seems to be no violation seeing as features are not highly correlated. Here we look into the one-hot encoded features.

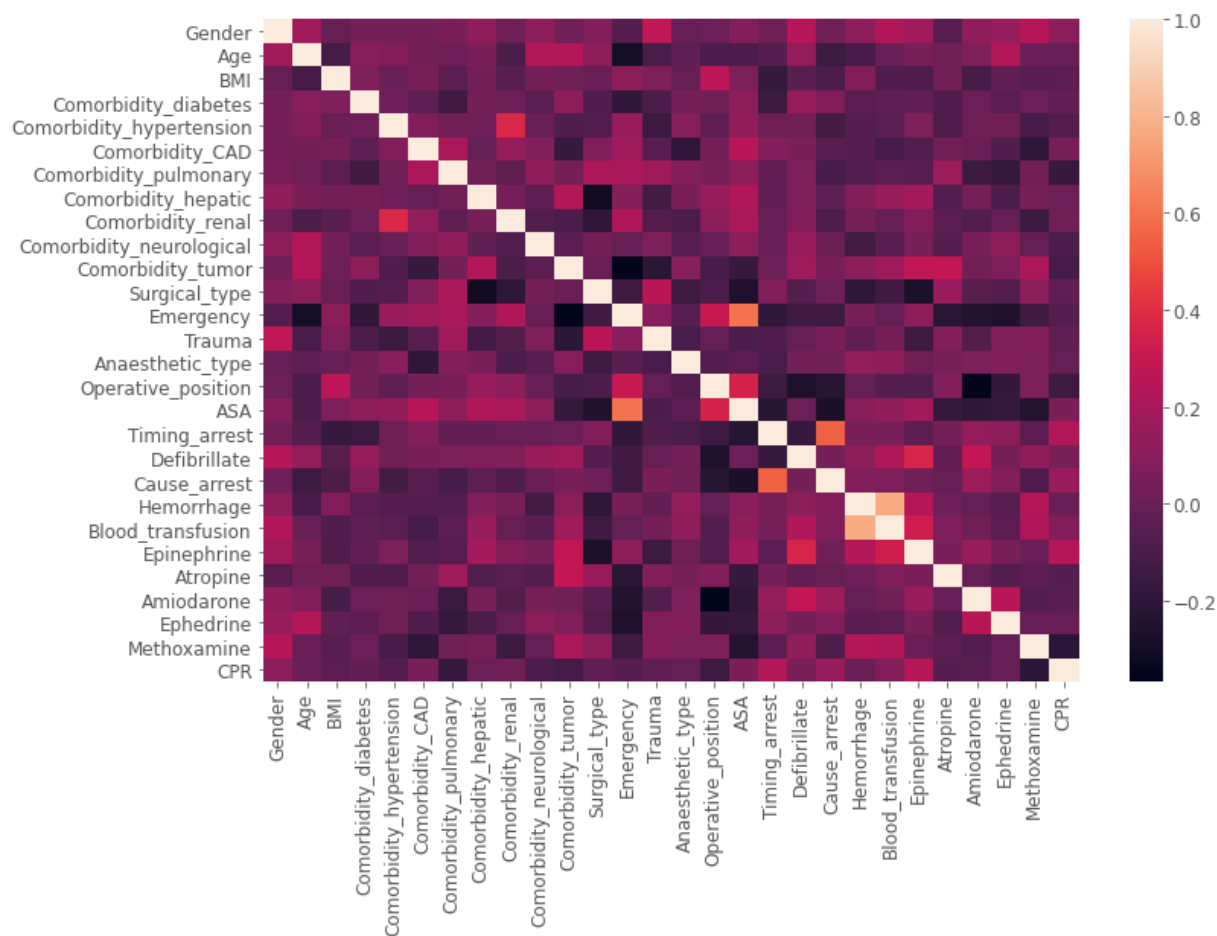
```
In [73]: 1 df_pdp.head()
```

Out[73]:

	Gender	Age	BMI	Comorbidity_diabetes	Comorbidity_hypertension	Comorbidity_CAD	Comorbi
0	1	2	1	0	1	0	
1	0	2	3	0	0	0	
2	0	1	1	0	0	1	
3	1	2	3	0	1	1	
4	1	2	2	0	1	1	

```
In [74]: 1 # calculate the correlation matrix
2 corr = df_pdp.corr()
3
4 # plot the heatmap
5 sns.heatmap(corr,
6             xticklabels=corr.columns,
7             yticklabels=corr.columns)
```

Out[74]: <AxesSubplot:>



Then I simply extract the features that have the highest absolute correlation by unstacking the correlation matrix and quicksorting it.


```
In [75]: 1 c = df_pdp.corr().abs()
2 s = c.unstack()
3 so = s.sort_values(kind="quicksort")
4 df_corr = pd.DataFrame(so).reset_index().dropna()
5 df_corr.columns = ['feature1', 'feature2', 'r']
6 df_corr = df_corr[df_corr.r < 1].sort_values('r', ascending=False)
7 df_corr.head(10)
```

Out[75]:

	feature1	feature2	r
755	Hemorrhage	Blood_transfusion	0.774969
754	Blood_transfusion	Hemorrhage	0.774969
753	ASA	Emergency	0.595407
752	Emergency	ASA	0.595407
751	Timing_arrest	Cause_arrest	0.542557
750	Cause_arrest	Timing_arrest	0.542557
749	Comorbidity_hypertension	Comorbidity_renal	0.378397
748	Comorbidity_renal	Comorbidity_hypertension	0.378397
747	Operative_position	Amiodarone	0.367689
746	Amiodarone	Operative_position	0.367689

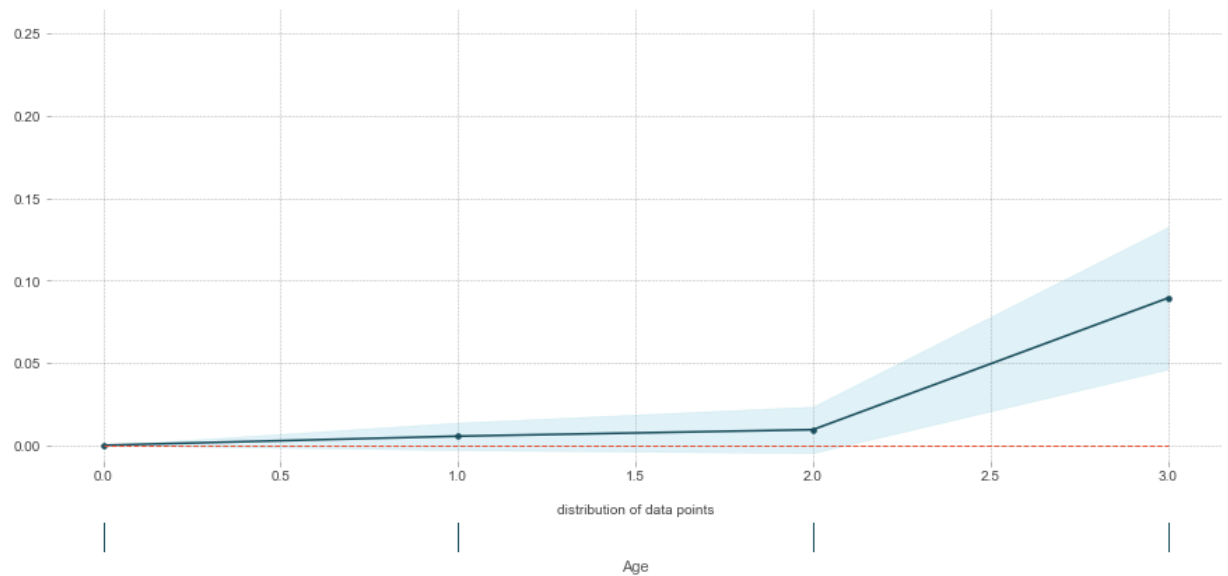
4.1 PDP - Single feature

The PDP plot for the feature "Age" shows that until the age of 50 there is a higher chance of earning more as a persons age increases. However, after the age of 50 we see this trend going the other direction, namely that age has a negative effect on the likelihood of earning more.

```
In [76]: 1 pdp_fare = pdp.pdp_isolate(  
2         model=rf_clf, dataset=df_pdp[df_pdp.columns[0:28]], model_features=df_pd  
3         )  
4 fig, axes = pdp.pdp_plot(pdp_fare, 'Age', plot_pts_dist=True)
```

PDP for feature "Age"

Number of unique grid points: 4



```
In [ ]: 1
```

