

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import plotly.graph_objs as go
6 from plotly.subplots import make_subplots
7 import os
8 from urllib.request import urlretrieve
9 from sklearn.model_selection import train_test_split
10 from imblearn.over_sampling import SMOTE
11

```

Il dataset preso in considerazione contiene informazioni utili per la predizione di infarti. Il dataset contiene le informazioni di 5000 pazienti, per ognuno di essi è specificato se hanno avuto un'infarto (stroke = 1) o no (stroke = 0), l'obiettivo è predire questa variabile discreta

```

1 def download(file, url):
2     if not os.path.isfile(file):
3         urlretrieve(url, file)

```

```

1 download("strokes.csv", "https://raw.githubusercontent.com/giacomoaccursi/proget

```

```

1 data = pd.read_csv("strokes.csv")

```

1. id: identificatore unico
2. gender: il genere della persone. "Male", "Female" or "Other"
3. age: l'età del paziente
4. hypertension: 0 se il paziente non soffre di ipertensione, 1 se ne soffre
5. heart_disease: 0 se il paziente non ha malattie cardiache, 1 se ne ha
6. ever_married: "No" se non si è mai sposato o "Yes" altrimenti
7. work_type: "children", "Govt_jov", "Never_worked", "Private" or "Self-employed"
8. Residence_type: "Rural" se vive in campagna o "Urban" se vive in città
9. avg_glucose_level: livello medio di glucosio nel sangue
10. bmi: indice di massa corporea, rapporto fra peso e quadrato dell'altezza
11. smoking_status: "formerly smoked" se ha fumato in precedenza, "never smoked" se non ha mai fumato, "smokes" se fuma attualmente o "Unknown" se il dato è sconosciuto
12. stroke: 1 se il paziente ha avuto un infarto, 0 altrimenti.

Essendo la variabile stroke discreta, siamo di fronte ad un problema di classificazione con due classi.

```

1 data.info()

<class 'pandas.core.frame.DataFrame'>

```

```

RangeIndex: 5110 entries, 0 to 5109
Data columns (total 12 columns):
#   Column                Non-Null Count  Dtype
---  -
0   id                     5110 non-null   int64
1   gender                 5110 non-null   object
2   age                    5110 non-null   float64
3   hypertension           5110 non-null   int64
4   heart_disease          5110 non-null   int64
5   ever_married           5110 non-null   object
6   work_type              5110 non-null   object
7   Residence_type         5110 non-null   object
8   avg_glucose_level      5110 non-null   float64
9   bmi                    4909 non-null   float64
10  smoking_status         5110 non-null   object
11  stroke                 5110 non-null   int64
dtypes: float64(3), int64(4), object(5)
memory usage: 479.2+ KB

```

La variabile id identifica univocamente il paziente. Non è quindi utile nella predizione. Può essere usata come indice della riga.

```
1 data.set_index('id', inplace=True)
```

```
1 data.head()
```

	gender	age	hypertension	heart_disease	ever_married	work_type	Residence_type
id							
9046	Male	67.0	0	1	Yes	Private	
51676	Female	61.0	0	0	Yes	Self-employed	
31112	Male	80.0	0	1	Yes	Private	
60182	Female	49.0	0	0	Yes	Private	
1665	Female	79.0	1	0	Yes	Self-employed	

Rinominiamo la colonna Residence_type cosicchè tutte le feature abbiano la prima lettera minuscola

```
1 data.rename(columns={'Residence_type': 'residence_type'}, inplace=True)
```

```
1 data['ever_married'] = data["ever_married"].replace({'No' : 0, 'Yes' : 1})
```

```
1 data.describe()
```

	age	hypertension	heart_disease	ever_married	avg_glucose_level
count	5110.000000	5110.000000	5110.000000	5110.000000	5110.000000
mean	43.226614	0.097456	0.054012	0.656164	106.147677
std	22.612647	0.296607	0.226063	0.475034	45.283560
min	0.080000	0.000000	0.000000	0.000000	55.120000
25%	25.000000	0.000000	0.000000	0.000000	77.245000
50%	45.000000	0.000000	0.000000	1.000000	91.885000
75%	61.000000	0.000000	0.000000	1.000000	114.090000
max	82.000000	1.000000	1.000000	1.000000	271.740000

- L'età media dei pazienti è circa 43 anni, il più giovane ha meno di un anno e il più anziano ne ha 42. Il 50% degli esaminati ha più di 45 anni.

```
1 (data.describe(exclude = ['float', 'int64']))
```

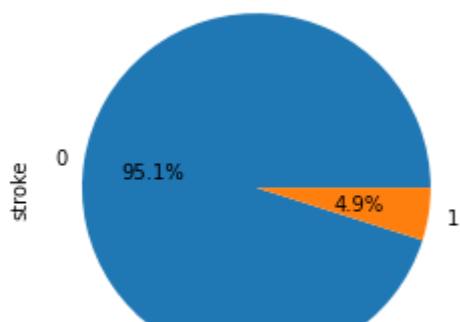
	gender	work_type	residence_type	smoking_status
count	5110	5110	5110	5110
unique	3	5	2	4
top	Female	Private	Urban	never smoked
freq	2994	2925	2596	1892

In un secondo momento occorrerà trattare i valori bmi nulli

```
1 data.isnull().sum()
```

```
gender          0
age             0
hypertension    0
heart_disease   0
ever_married    0
work_type       0
residence_type  0
avg_glucose_level 0
bmi            201
smoking_status  0
stroke          0
dtype: int64
```

```
1 data['stroke'].value_counts().plot.pie(autopct='%1.1f%%');
```



Si noti come il dataset risulti fortemente sbilanciato: il numero di persone che hanno avuto un infarto è nettamente inferiore a quelle che non l'hanno avuto. Nelle fasi successive si cercherà di risolvere questo problema

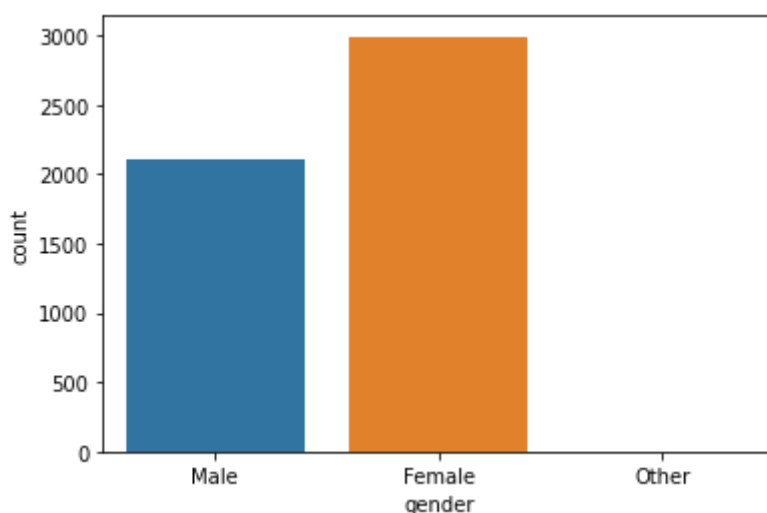
▼ Data Exploration

Circa il 60% delle persone in esame sono donne.

```
1 sns.countplot(data['gender']);
```

/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning

Pass the following variable as a keyword arg: x. From version 0.12, the only



Sembra che il sesso non influisca sulla probabilità di avere un infarto.

```
1 fig = plt.figure(figsize=(20, 10))
2 plt.subplot(2, 3, 1)
3 (data[data["stroke"] == 1]["gender"]).value_counts().plot.pie(autopct='%1.1f%%')
4 plt.subplot(2, 3, 2)
5 sns.countplot((data[data["gender"] == "Female"]["stroke"]))
6 plt.subplot(2, 3, 3)
7 sns.countplot((data[data["gender"] == "Male"]["stroke"]))
```

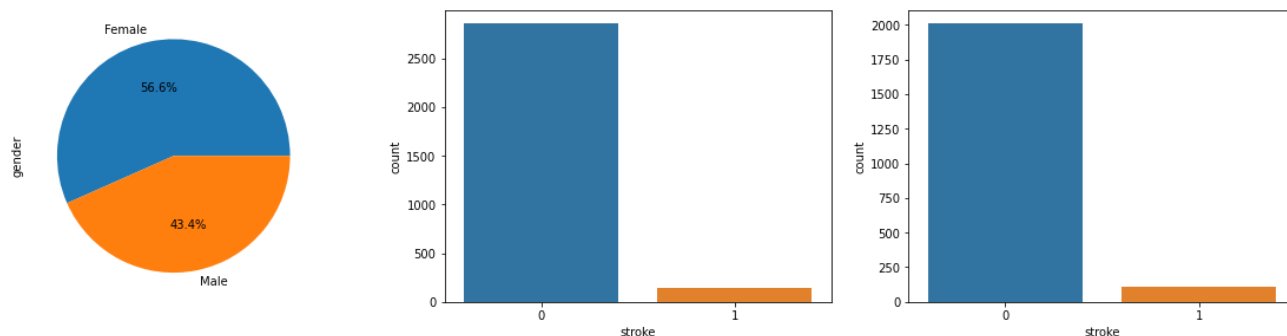
```
7 sns.countplot(data[data['gender'] == 'male']['stroke'])
8 fig.show()
```

/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning

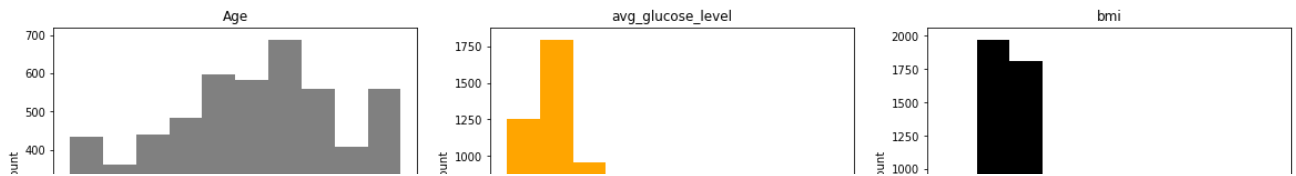
Pass the following variable as a keyword arg: x. From version 0.12, the only

/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning

Pass the following variable as a keyword arg: x. From version 0.12, the only



```
1 plt.figure(figsize=(20, 10))
2
3 plt.subplot(2, 3, 1)
4 plt.title('Age')
5 plt.hist(data['age'], label="weight", color='gray')
6 plt.ylabel('count')
7 plt.xlabel('Year')
8
9
10 plt.subplot(2, 3, 2)
11 plt.title('avg_glucose_level')
12 plt.hist(data['avg_glucose_level'], label="weight", color="orange")
13 plt.ylabel('count')
14 plt.xlabel('mg/dl')
15
16 plt.subplot(2, 3, 3)
17 plt.title('bmi')
18 plt.hist(data['bmi'], label="weight", color="black")
19 plt.ylabel('count')
20 plt.xlabel('kg/m2')
21
22 plt.show()
```



```
1 ((data["bmi"] > 25) & (data["bmi"] < 30)).mean()
```

```
0.27045009784735813
```

Year

mg/dl

kg/m2

```
1 (data["bmi"] > 30).mean()
```

```
0.3704500978473581
```

Non è specificato se i valori di glucosio nel sangue siano stati rilevati a digiuno o dopo 120' dal carico orale di glucosio. Secondo l'indicatore bmi il 27% delle persone in esame risulta sovrappesa e il 37% obesa, si noti però che il bmi non è un buon indicatore per il fatto che tiene in considerazione solo l'altezza in relazione al peso senza considerare la percentuale di grasso corporeo.

Valori glicemici a digiuno:

- Normale : 80-99 mg/dl
- Alterata : 100-125 mg/dl
- Diabete : >126

Valori glicemici dopo 120' dal carico orale di glucosio:

- Normale : <140 mg/dl
- Alterata : 140-200 mg/dl
- Diabete : >200

valori bmi:

- Sottopeso : <18.5
- Normale : 18.5 - 24.9
- Sovrappeso : 25- 29.9
- Obeso : >30

Si noti come per valori bmi superiori a 30 la probabilità di avere infarti aumenti leggermente e come livelli di glucosio superiori a 150 corrispondano ad un numero maggiore di infarti.

```
1 plt.figure(figsize=(20, 10))
2
3 plt.subplot(2, 3, 1)
4 sns.distplot(data[data['stroke'] == 0]["bmi"], color='green') # No Stroke - green
5 sns.distplot(data[data['stroke'] == 1]["bmi"], color='red') # Stroke - Red
6
7 plt.title('No Stroke vs Stroke by BMI', fontsize=15)
```

```

8 plt.xlim([10, 100])
9
10 plt.subplot(2, 3, 2)
11 sns.distplot(data[data['stroke'] == 0]["avg_glucose_level"], color='green') # No
12 sns.distplot(data[data['stroke'] == 1]["avg_glucose_level"], color='red') # Str
13
14 plt.title('No Stroke vs Stroke by Avg. Glucose Level', fontsize=15)
15 plt.xlim([30, 330])
16 plt.show()

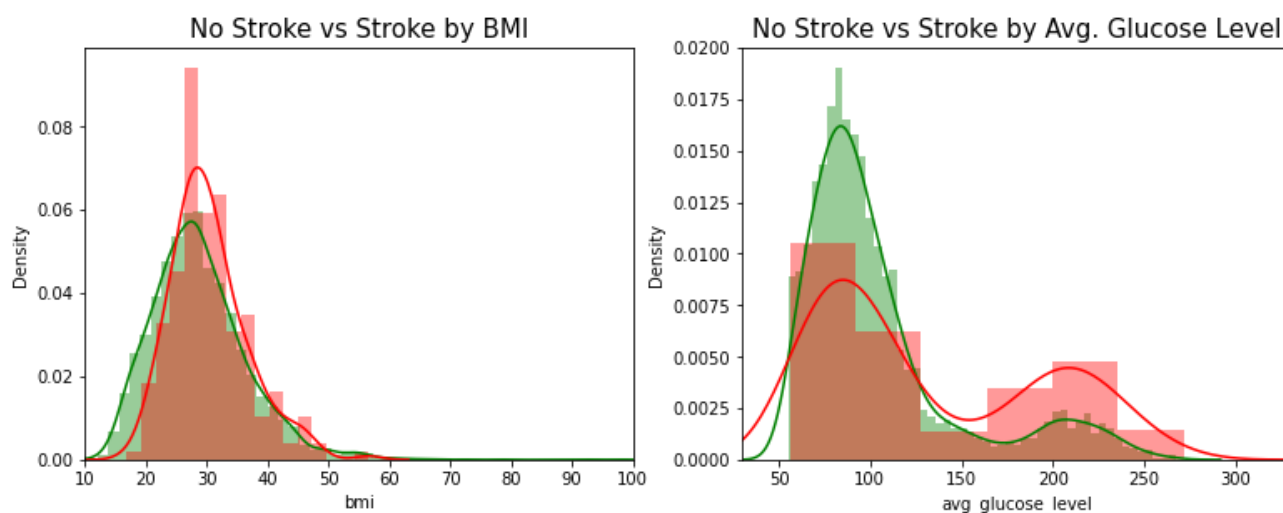
```

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2557: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Use `displot` instead.

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2557: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Use `displot` instead.

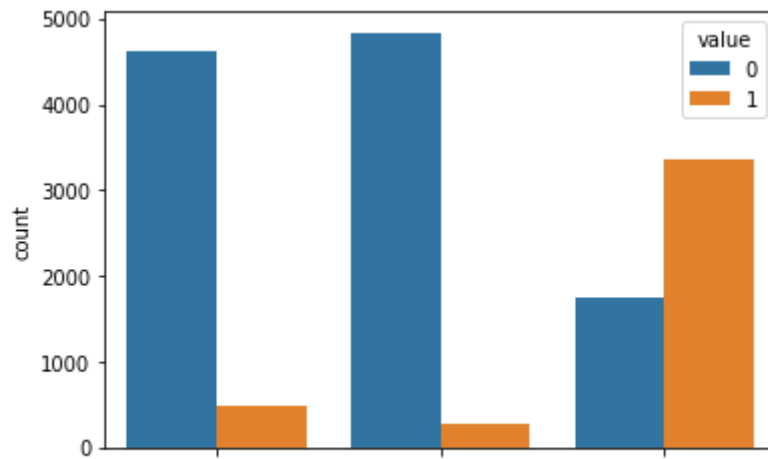
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2557: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Use `displot` instead.

/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2557: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Use `displot` instead.



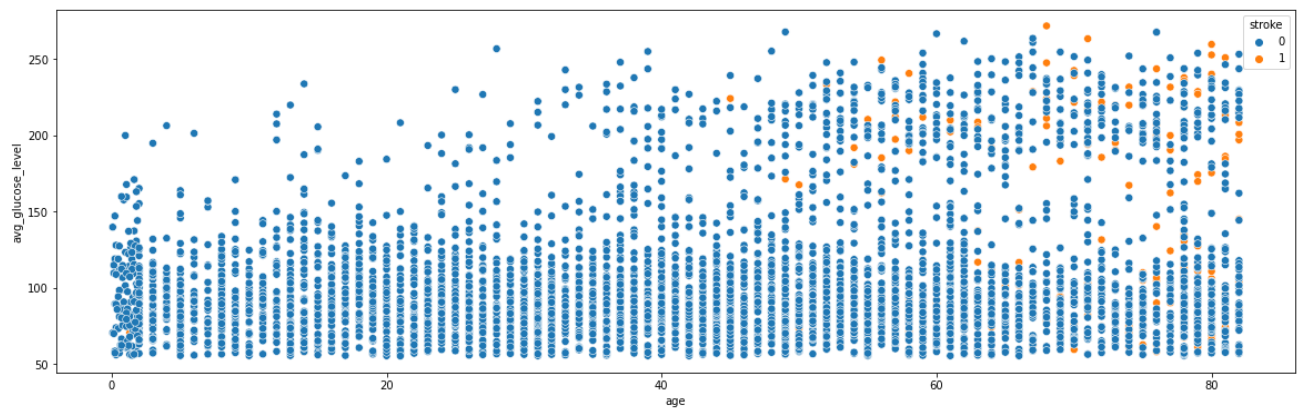
Dal seguente istogramma emerge come le persone che sono state o sono sposate abbiamo una probabilità maggiore di avere un infarto, tuttavia questo è molto probabilmente causato dal fatto che le persone sposate o che sono state sposate hanno un'età maggiore delle persone che non si sono sposate.

```
1 sns.countplot(x="variable", hue="value", data= pd.melt(data.loc[:, ['hypertensic
```



Dal seguente scatterplot si nota come gli infarti sono più diffusi nelle persone in età avanzata e con livelli di glucosio alti.

```
1 fig = plt.gcf()
2 fig.set_size_inches(20, 6)
3 sns.scatterplot(x=data['age'], y=data['avg_glucose_level'], hue=data['stroke'], s=
```

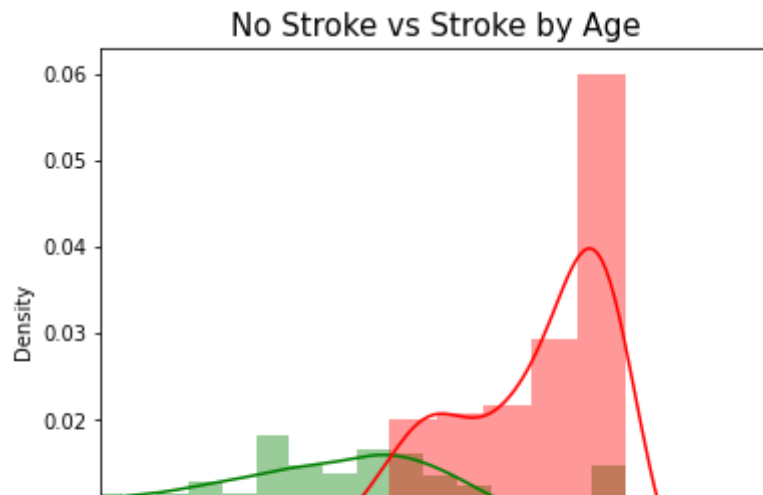


Come previsto il numero di infarti è maggiore nelle persone più anziane.

```
1 plt.figure(figsize=(6,5))
2
3 sns.distplot(data[data['stroke'] == 0]["age"], color='green') # No Stroke - green
4 sns.distplot(data[data['stroke'] == 1]["age"], color='red') # Stroke - Red
5
6 plt.title('No Stroke vs Stroke by Age', fontsize=15)
7 plt.xlim([18,100])
8 plt.show()
```

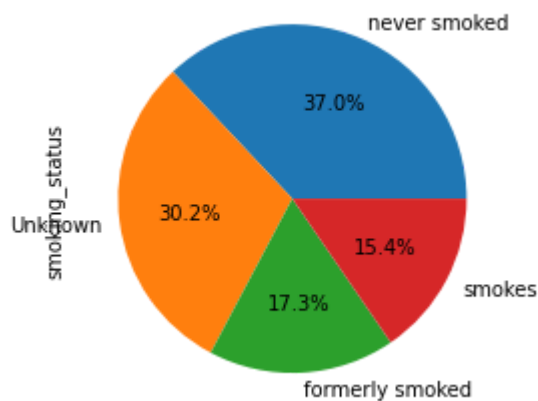


```
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2557: FutureWarning:
`distplot` is a deprecated function and will be removed in a future version. Use
`displot` instead.
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2557: FutureWarning:
`distplot` is a deprecated function and will be removed in a future version. Use
```



Il numero di infarti è maggiore nelle persone che non hanno mai fumato rispetto a quelle che hanno fumato o fumano tuttora. Ricordiamoci comunque che del 30% del dataset in questione non abbiamo informazioni.

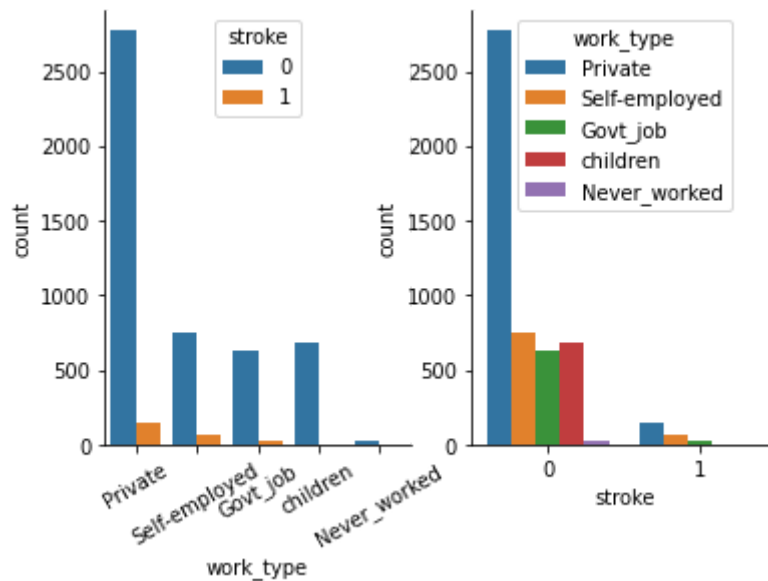
```
1 data['smoking_status'].value_counts().plot.pie(autopct='%1.1f%%');
```



Non sembra che ci siano lavori che influenzano la probabilità di avere un infarto

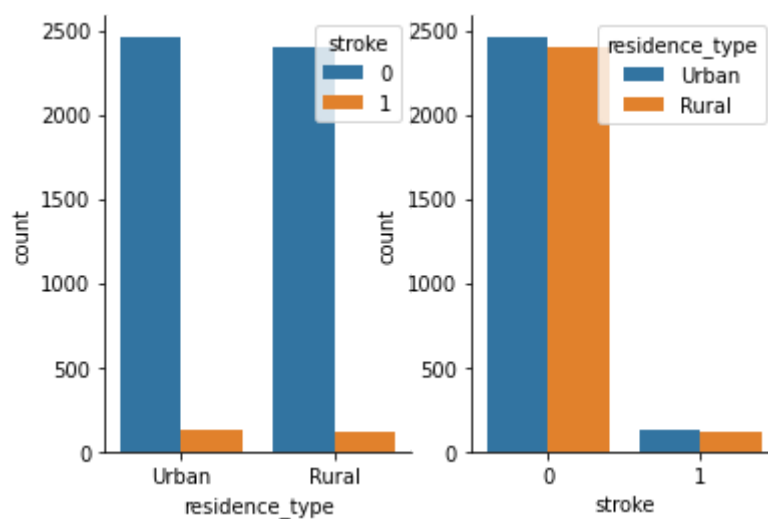
```
1 plt.subplot(1, 2, 1)
2
3 sns.countplot(x='work_type', hue='stroke', data=data)
4 plt.xticks(rotation=30)
5 sns.despine()
6
7 plt.subplot(1, 2, 2)
8 sns.countplot(x='stroke', hue='work_type', data=data)
9 sns.despine()
```

```
10
11 plt.show()
```



Lo stesso discorso vale per la residenza delle persone: vivere in città o in campagna non sembra rilevante

```
1 plt.subplot(1, 2, 1)
2
3 sns.countplot(x='residence_type', hue='stroke', data=data)
4 sns.despine()
5
6 plt.subplot(1, 2, 2)
7 sns.countplot(x='stroke', hue='residence_type', data=data)
8 sns.despine()
9
10 plt.show()
```



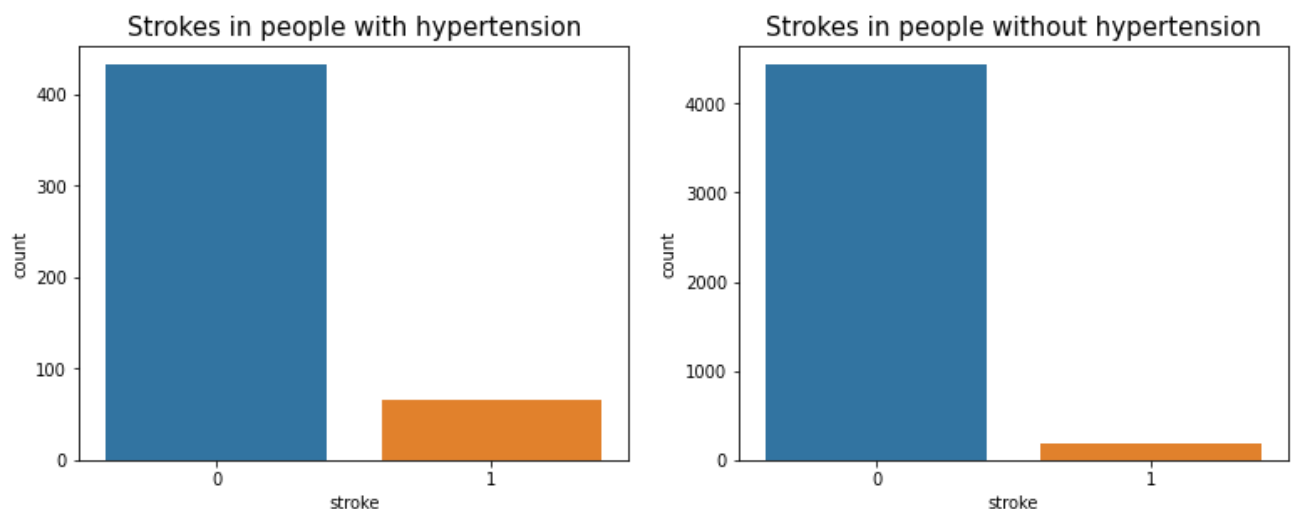
I seguenti istogrammi evidenziano come problemi al cuore e ipertensione aumentino la probabilità di avere un infarto

```

1 plt.figure(figsize=(20, 10))
2 plt.subplot(2, 3, 1)
3 plt.title('Strokes in people with hypertension', fontsize=15)
4 sns.countplot((data[data["hypertension"] == 1]["stroke"]))
5 plt.subplot(2, 3, 2)
6 plt.title('Strokes in people without hypertension', fontsize=15)
7 sns.countplot((data[data["hypertension"] == 0]["stroke"]))
8 fig.show()

```

/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning
Pass the following variable as a keyword arg: x. From version 0.12, the only
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning
Pass the following variable as a keyword arg: x. From version 0.12, the only



```

1 plt.figure(figsize=(20, 10))
2 plt.subplot(2, 3, 1)
3 plt.title('Strokes in people with hearth disease', fontsize=15)
4 sns.countplot((data[data["heart_disease"] == 1]["stroke"]))
5 plt.subplot(2, 3, 2)
6 plt.title('Strokes in people without hearth disease', fontsize=15)
7 sns.countplot((data[data["heart_disease"] == 0]["stroke"]))
8 fig.show()

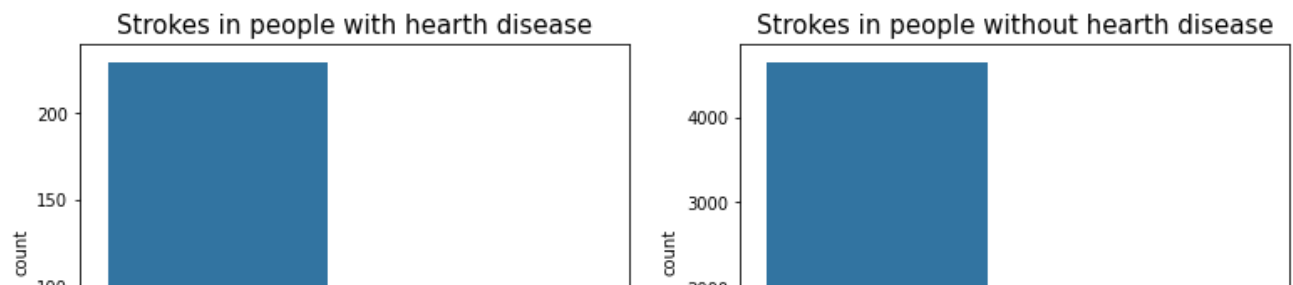
```

/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning

Pass the following variable as a keyword arg: x. From version 0.12, the only

/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning

Pass the following variable as a keyword arg: x. From version 0.12, the only



```
1 plt.figure(figsize=(15, 15))
2 plt.subplot(4,2,1)
3 sns.countplot(data['work_type'])
4 plt.subplot(4,2,2)
5 sns.countplot(data['residence_type'])
6 plt.subplot(4,2,3)
7 sns.countplot(data['smoking_status'])
8 plt.subplot(4,2,4)
9 sns.countplot(data['ever_married'])
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning
```

```
Pass the following variable as a keyword arg: x. From version 0.12, the only
```

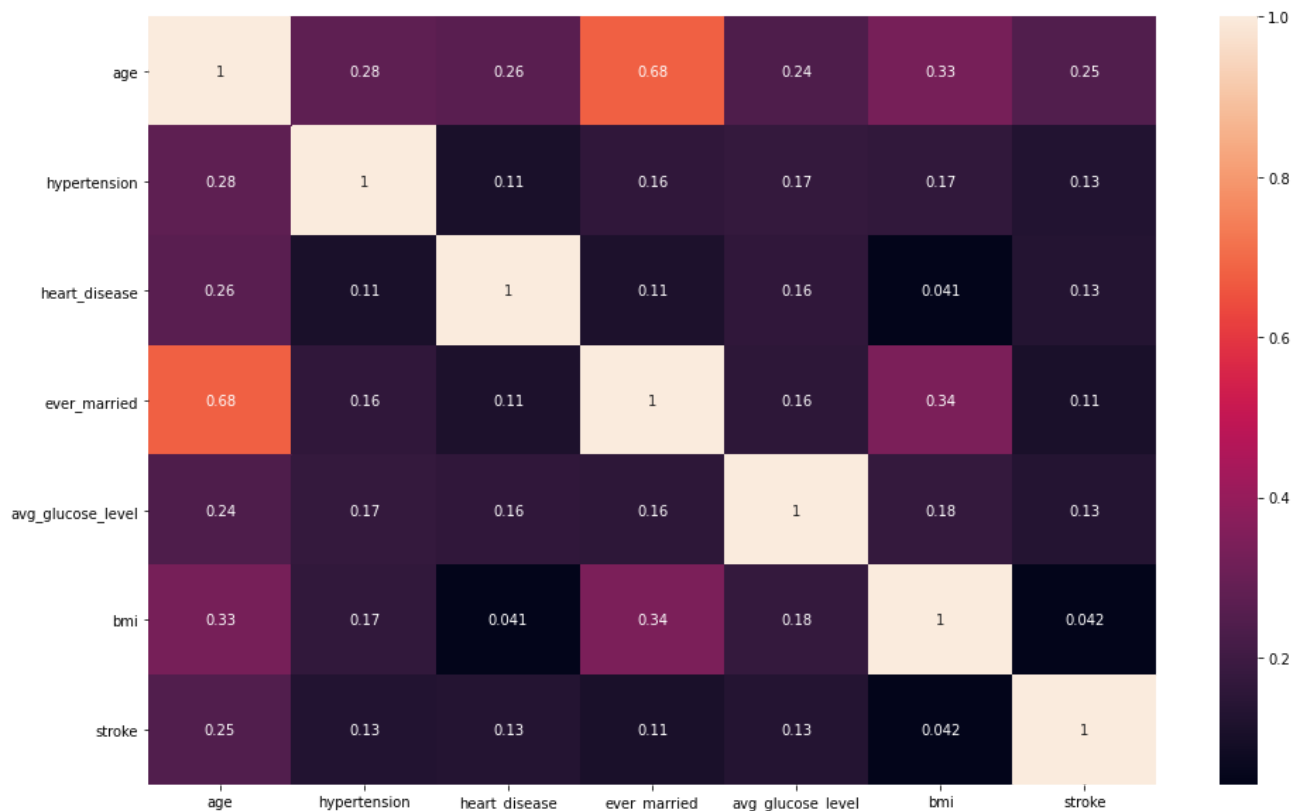
```
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning
```

```
Pass the following variable as a keyword arg: x. From version 0.12, the only
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning
```

```
Pass the following variable as a keyword arg: x. From version 0.12, the only
```

```
1 plt.figure(figsize=(16,10))
2 sns.heatmap(data.corr(method='pearson'), annot=True);
```



Come ipotizzato in precedenza esiste una correlazione abbastanza forte fra l'età e il numero di persone sposate.

▼ Feature engineering

```
1 data.isnull().sum()
```

```
gender      0
```

```

age                0
hypertension       0
heart_disease      0
ever_married       0
work_type          0
residence_type     0
avg_glucose_level  0
bmi                201
smoking_status     0
stroke             0
dtype: int64

```

Eliminiamo i record con ibm nullo

```
1 data.dropna(inplace=True)
```

```
1 data.isnull().sum()
```

```

gender            0
age               0
hypertension      0
heart_disease     0
ever_married      0
work_type         0
residence_type    0
avg_glucose_level 0
bmi               0
smoking_status    0
stroke            0
dtype: int64

```

```
1 data["age"] = data["age"].astype(int)
```

```
2 data["avg_glucose_level"] = data["avg_glucose_level"].astype(int)
```

```
3 data["bmi"] = data["bmi"].astype(int)
```

Il dataset risulta fortemente sbilanciato, creiamo quindi nuovi record. E' stato utilizzato SMOTENC al posto di SMOTE poichè è più indicato nei casi in cui si hanno sia variabili categoriche che numeriche

```
1 categorical_features = ["gender", "work_type", "residence_type", "smoking_status"]
```

```
2 data = pd.get_dummies(data, columns=categorical_features, prefix=categorical_features)
```

```
1 data
```

	age	hypertension	heart_disease	ever_married	avg_glucose_level	bmi
id						
9046	67	0	1	1	228	36
31112	80	0	1	1	105	32
60182	49	0	0	1	171	34
1665	79	1	0	1	174	24
56669	81	0	0	1	186	29
...
14180	13	0	0	0	103	18
44873	81	0	0	1	125	40

```

1 sm = SMOTE(random_state=42)
2 y = data["stroke"]
3 X = data.drop("stroke", axis=1)
4 X, y = sm.fit_resample(X, y)
5

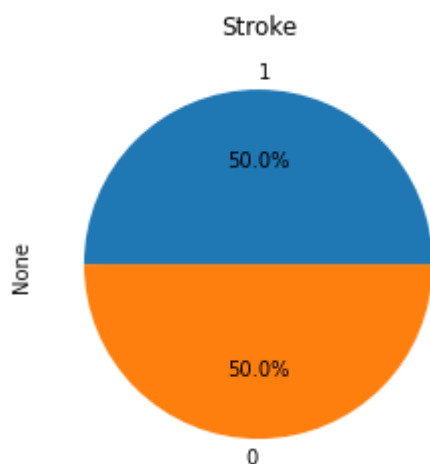
```

```

/usr/local/lib/python3.7/dist-packages/sklearn/externals/six.py:31: FutureWarning:
The module is deprecated in version 0.21 and will be removed in version 0.23 :
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:144: FutureWarning:
The sklearn.neighbors.base module is deprecated in version 0.22 and will be removed in version 0.24 :
/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning:
Function safe_indexing is deprecated; safe_indexing is deprecated in version 0.21 and will be removed in version 0.23 :

```

```
1 pd.value_counts(y).plot.pie(autopct="%.1f%", title="Stroke");
```



```
1 from sklearn.compose import make_column_transformer
```

```

1 from sklearn.compose import make_column_transformer
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.model_selection import train_test_split
4 from sklearn.model_selection import KFold
5 from sklearn.pipeline import Pipeline
6 from sklearn.metrics import classification_report
7 from sklearn.metrics import mean_squared_error
8 from sklearn.metrics import confusion_matrix
9 from sklearn.model_selection import GridSearchCV
10 from sklearn.preprocessing import PolynomialFeatures
11 from sklearn.metrics import precision_score, recall_score, f1_score
12

```

```

1 scaler = StandardScaler()
2 X = scaler.fit_transform(X)

```

```

1 X_train, X_val, y_train, y_val = train_test_split(
2     X,y,
3     test_size = 1/3,
4     random_state = 42
5 )

```

```

1 models = {}
2 kf = KFold(n_splits=5, shuffle=True, random_state=42)

```

```

1 def accuracy_interval(f):
2     N = len(y_val)
3     n_min = f + ( 1.96**2/(2*N) - 1.96 * np.sqrt( (f/N) - (f**2/N) + (1.96**2/(4*N)
4     n_max = f + ( 1.96**2/(2*N) + 1.96 * np.sqrt( (f/N) - (f**2/N) + (1.96**2/(4*N)
5     d = 1 + (1.96**2 / N)
6     e_min = n_min / d
7     e_max = n_max / d
8     return np.round(e_min,4), np.round(e_max,4)

```

```

1 from sklearn.linear_model import Perceptron
2
3 std_perceptron = Perceptron(n_jobs=-1, early_stopping=True, n_iter_no_change=5)
4
5 parameters = {
6     'penalty': [None, 'l1', 'l2', 'elasticnet'],
7     'alpha': [0.0001, 0.001, 0.01, 1],
8     'tol': [1e-9, 1e-6, 1e-3, 1, 1e3, 1e6],
9 }
10
11 perceptron_cv = GridSearchCV(std_perceptron, parameters, cv=kf, n_jobs=-1)
12 perceptron_cv.fit(X_train, y_train)
13
14 print('Best parameters:', perceptron_cv.best_params_)
15
16 models["Perceptron"] = {"Model": perceptron_cv.best_estimator_ , "Score" : perce

```

Best parameters: {'alpha': 0.0001, 'penalty': 'l1', 'tol': 1e-09}


```

1 poly_perceptron = Pipeline([
2     ('poly', PolynomialFeatures(degree=3)),
3     ('perceptron', Perceptron(n_jobs=-1, early_stopping=True, n_iter_no_change=5
4 ]))
5
6 parameters = {
7     'perceptron__penalty': ['l1', 'l2'],
8     'perceptron__alpha': [0.0001, 0.001, 0.01],
9     'perceptron__tol': [1e-9, 1e-6, 1e-3, 1],
10 }
11
12 poly_perceptron_cv = GridSearchCV(poly_perceptron, parameters, cv=kf, n_jobs=-1,
13 poly_perceptron_cv.fit(X_train, y_train)
14
15 print('Best parameters:', poly_perceptron_cv.best_params_)
16
17 models["Perceptron"] = {"Model": perceptron_cv.best_estimator_ , "Score" : poly_

```

Best parameters: {'perceptron__alpha': 0.0001, 'perceptron__penalty': 'l1', 'l1'

```

1 from sklearn.linear_model import LogisticRegression
2
3 std_lr = LogisticRegression(solver="saga", random_state=42)
4 grid = [
5     {
6         "penalty": ["l2", "l1"],
7         "C": [0.1, 1, 10],
8         "tol" : [1e-9, 1e-6, 1e-3, 1e-2, 1e-1, 1]
9     },
10    {
11        "penalty": ["elasticnet"],
12        "C": [0.1, 1, 10],
13        "l1_ratio": [0.2, 0.5],
14        "tol" : [1e-9, 1e-6, 1e-3, 1e-2, 1e-1, 1]
15    }
16 ]
17
18 lr_gs = GridSearchCV(std_lr, grid, cv=kf, n_jobs=-1, return_train_score=True)
19 lr_gs.fit(X_train, y_train)
20
21 print('Best parameters:', lr_gs.best_params_)
22
23 models["Logistic Regression"] = {"Model": lr_gs.best_estimator_ , "Score" : lr_g

```

Best parameters: {'C': 0.1, 'l1_ratio': 0.5, 'penalty': 'elasticnet', 'tol': (

```

1 from sklearn.neighbors import KNeighborsClassifier
2 knc = KNeighborsClassifier(n_jobs=-1)
3 grid = {
4     'n_neighbors': range(1, 10, 1),
5     'weights': ['uniform', 'distance']

```

```

6 }
7
8 knc_gs = GridSearchCV(knc, grid, cv=kf, n_jobs=-1, return_train_score=True)
9 knc_gs.fit(X_train, y_train)
10
11 print('Best parameters:', knc_gs.best_params_)
12
13 models["KNeighborsClassifier"] = {"Model": knc_gs.best_estimator_ , "Score" : kr

```

Best parameters: {'n_neighbors': 2, 'weights': 'uniform'}

```

1 from sklearn.svm import SVC
2 std_svm = SVC()
3
4 grid = {
5     'kernel': ['rbf'],
6     'C': [0.01, 0.1, 1],
7 }
8
9 svm_gs = GridSearchCV(std_svm, grid, cv=kf, n_jobs=-1, return_train_score=True)
10 svm_gs.fit(X_train, y_train)
11
12 print('Best parameters:', lr_gs.best_params_)
13
14 models["Support Vector Machine"] = {"Model": svm_gs.best_estimator_ , "Score" :

```

Best parameters: {'C': 0.1, 'l1_ratio': 0.5, 'penalty': 'elasticnet', 'tol': (

```

1 from sklearn.tree import DecisionTreeClassifier
2
3 d_tree = DecisionTreeClassifier(random_state=42)
4
5 grid = {
6     'min_samples_split': range(2, 4, 1),
7     'min_samples_leaf': range(1, 4, 1),
8     'max_depth': [None] + [i for i in range(2, 7)],
9     'max_features': range(2, data.columns.size, 1)}
10
11
12 d_tree = GridSearchCV(d_tree, grid, cv=kf, n_jobs=-1, return_train_score=True)
13 d_tree.fit(X_train, y_train)
14
15 print('Best parameters:', d_tree.best_params_)
16
17 models["Decision Tree"] = {"Model": d_tree.best_estimator_ , "Score" : d_tree.sc

```

Best parameters: {'max_depth': None, 'max_features': 10, 'min_samples_leaf': :

```

1 from sklearn.ensemble import RandomForestClassifier
2
3 rfc = RandomForestClassifier(n_jobs=-1, random_state=3)
4

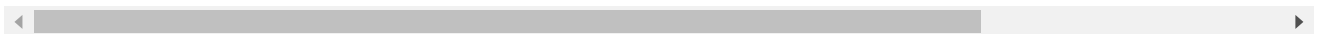
```

```

5 grid = {
6     'n_estimators': [100, 200, 300],
7     'max_depth': [2, 4, 6, 8, 10],
8     'min_samples_leaf': [1, 2, 4],
9     'min_samples_split': [2, 5, 10],
10 }
11
12
13 rfc_gs = GridSearchCV(rfc, grid, cv=kf, n_jobs=-1, return_train_score=True)
14 rfc_gs.fit(X_train, y_train)
15
16 print('Best parameters:', rfc_gs.best_params_)
17
18 models["Random Forest"] = {"Model": rfc_gs.best_estimator_ , "Score" : rfc_gs.sc

```

Best parameters: {'max_depth': 10, 'min_samples_leaf': 1, 'min_samples_split':

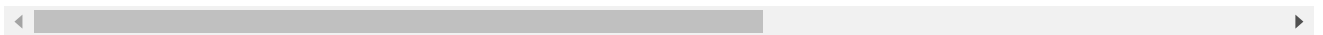


```

1 from xgboost import XGBClassifier
2
3 std_xgb = XGBClassifier(nthread=8, objective='binary:logistic')
4
5 parameters = {
6     'xgb__eta': [0.002, 0.1, 0.5],
7     'xgb__min_child_weight': [4, 10],
8     'xgb__max_depth': [6],
9     'xgb__n_estimators': [150, 300],
10    'xgb__alpha': [0.0001, 0.001]
11 }
12
13 xgb_gs = GridSearchCV(std_xgb, parameters, cv=kf, n_jobs=-1, return_train_score=
14 xgb_gs.fit(X_train, y_train)
15
16 print('Best parameters:', xgb_gs.best_params_)
17
18 models["XGBoost"] = {"Model": xgb_gs.best_estimator_ , "Score" : xgb_gs.score}

```

Best parameters: {'xgb__alpha': 0.0001, 'xgb__eta': 0.002, 'xgb__max_depth': 6



```

1 from sklearn.dummy import DummyClassifier
2 random = DummyClassifier(strategy="uniform", random_state=42)
3 random.fit(X_train, y_train)
4
5 models["Random"] = {"Model": random , "Score" : random.score}

```

▼ Model Comparison

```

1 for name, model in models.items():
2     model = models[name]
3     y_pred = model["Model"].predict(X_val)
4     ...

```

```

4 print(name)
5 model["Precision"] = precision_score(y_val, y_pred, pos_label=0)
6 model["Recall"] = recall_score(y_val, y_pred, pos_label=0)
7 model["F1_Score"] = f1_score(y_val, y_pred, average="macro")
8 model["mse"] = mean_squared_error(y_val, y_pred)
9 print("Precision = ", model["Precision"])
10 print("Recall = ", model["Recall"])
11 print("F1 score = ", model["F1_Score"])
12 print(pd.DataFrame(confusion_matrix(y_val, y_pred), index=["No Stroke", "Stroke"], columns=["No Stroke", "Stroke"]))
13 print("Accuracy interval = ", accuracy_interval(model["F1_Score"]))
14 print("MSE = ", model["mse"])
15 print("\n\n")

```

```

Support Vector Machine
Precision = 0.9262295081967213
Recall = 0.9968494013862634
F1 score = 0.9580897781602455
      No Stroke  Stroke
No Stroke      1582      5
Stroke          126     1421
Accuracy interval = (0.9505, 0.9646)
MSE = 0.041799617102744095

```

```

Decision Tree
Precision = 0.9458128078817734
Recall = 0.9678638941398866
F1 score = 0.9556209452154114
      No Stroke  Stroke
No Stroke      1536      51
Stroke          88     1459
Accuracy interval = (0.9478, 0.9623)
MSE = 0.04435226547543076

```

```

Random Forest
Precision = 0.9451476793248945
Recall = 0.9880277252678009
F1 score = 0.9648562013706836
      No Stroke  Stroke
No Stroke      1568      19
Stroke          91     1456
Accuracy interval = (0.9578, 0.9708)
MSE = 0.03509891512444161

```

```

XGBoost
Precision = 0.9303423848878394
Recall = 0.9930686830497795
F1 score = 0.9587477871956984
      No Stroke  Stroke
No Stroke      1576      11
Stroke          118     1429
Accuracy interval = (0.9512, 0.9652)
MSE = 0.04116145500957243

```

```

Random
Precision = 0.5091714104996837
Recall = 0.5072463768115942
F1 score = 0.5028132127728902
      No Stroke  Stroke
No Stroke      805    782
Stroke         776    771
Accuracy interval = (0.4853, 0.5203)
MSE = 0.4971282705807275

```

Confrontiamo con una confidenza del 95% i modelli, dati i rispettivi errori.

```

1 from itertools import combinations
2
3 def intervall95(mse1, mse2, confidence):
4     z = 1.96
5     d = np.abs(mse1 - mse2)
6     variance = (mse1 * (1 - mse1)) / len(X_val) + (mse2 * (1 - mse2)) / len(X_val)
7     d_min = d - z * np.sqrt(variance)
8     d_max = d + z * np.sqrt(variance)
9     return d_min, d_max
10
11 svm_error = models["Support Vector Machine"]["mse"]
12 lre_error = models["Logistic Regression"]["mse"]
13 knc_error = models["KNeighborsClassifier"]["mse"]
14 tree_error = models["Decision Tree"]["mse"]
15 forest_error = models["Random Forest"]["mse"]
16 xgb_error = models["XGBoost"]["mse"]
17
18 mse = [(svm_error, "svm"), (lre_error, "Logistic Regression"),
19        (knc_error, "knc"), (tree_error, "Decision Tree"), (xgb_error, "XGBoost")]
20
21 print (f"{'Models':<40} {'Interval':<15}")
22 for m1, m2 in list(combinations(mse, 2)):
23     mse1, mse2 = m1[0], m2[0]
24     name1, name2 = m1[1], m2[1]
25     comparison = name1 + " vs " + name2
26     print (f"{'comparison':<40} {np.round(intervall95(mse1, mse2, 0.95), 4)} ")

```

Models	Interval
svm vs Logistic Regression	[-0.0075 0.0126]
svm vs knc	[0.0006 0.0192]
svm vs Decision Tree	[-0.0075 0.0126]
svm vs XGBoost	[-0.0092 0.0105]
svm vs random forest	[-0.0028 0.0162]
Logistic Regression vs knc	[0.003 0.0219]
Logistic Regression vs Decision Tree	[-0.0102 0.0102]
Logistic Regression vs XGBoost	[-0.0068 0.0132]
Logistic Regression vs random forest	[-0.0004 0.0189]
knc vs Decision Tree	[0.003 0.0219]
knc vs XGBoost	[-0.0004 0.0185]
knc vs random forest	[-0.0057 0.0121]
Decision Tree vs XGBoost	[-0.0068 0.0132]

Decision Tree vs random forest	[-0.0004 0.0189]
XGBoost vs random forest	[-0.0034 0.0155]

```
1 for name, model in models.items():
2     print(f"{name:<30}{model['F1_Score']}")
```

Perceptron	0.9459809205597873
Logistic Regression	0.9555274334745528
KNeighborsClassifier	0.9680567191987883
Support Vector Machine	0.9580897781602455
Decision Tree	0.9556209452154114
Random Forest	0.9648562013706836
XGBoost	0.9587477871956984
Random	0.5028132127728902

Guardando l'F1 Score dei modelli si può notare come i punteggi più alti siano ottenuti da:

1. KNeighborsClassifier
2. Random Forest
3. XGBoost

tuttavia, dal confronto fra i modelli possiamo notare che non esistono differenze significative fra i modelli.

Si verifica che la differenza fra i modelli e uno casuale sia statisticamente significativa

```
1 def interval99(mse1, mse2, confidence):
2     z = 2.58
3     d = np.abs(mse1 - mse2)
4     variance = (mse1 * (1 - mse1)) / len(X_val) + (mse2 * (1 - mse2)) / len(X_val)
5     d_min = d - z * np.sqrt(variance)
6     d_max = d + z * np.sqrt(variance)
7     return d_min, d_max
8
9 mse_random = models["Random"]["mse"]
10 print (f"{'Models':<40} Interval")
11 for m in mse:
12     mse_i = m[0]
13     name_i = m[1]
14     comparison = name_i + " vs Random"
15     print (f"{'comparison':<40} {np.round(interval99(mse_i , mse_random, 0.99), 4)}")
```

Models	Interval
svm vs Random	[0.4305 0.4801]
Logistic Regression vs Random	[0.4279 0.4777]
knc vs Random	[0.4408 0.4896]
Decision Tree vs Random	[0.4279 0.4777]
XGBoost vs Random	[0.4312 0.4808]
random forest vs Random	[0.4375 0.4866]

✓ 0s completed at 3:02 PM ● ✕