

- **Objective of the analysis**

The objective of this analysis is to be able to define a model to predict if a patient can have an heart attack based on some heart parameters

- **Description of the data set**

For this project I'm going to use a Data Set representing Heart Attack classification. Each case of heart attack has some parameters:

1. age - age in years
2. sex - sex (1 = male; 0 = female)
3. cp - chest pain type (1 = typical angina; 2 = atypical angina; 3 = non-anginal pain; 0 = asymptomatic)
4. trestbps - resting blood pressure (in mm Hg on admission to the hospital)
5. chol - serum cholesterol in mg/dl
6. fbs - fasting blood sugar > 120 mg/dl (1 = true; 0 = false)
7. restecg - resting electrocardiographic results (1 = normal; 2 = having ST-T wave abnormality; 0 = hypertrophy)
8. thalach - maximum heart rate achieved
9. exang - exercise induced angina (1 = yes; 0 = no)
10. oldpeak - ST depression induced by exercise relative to rest
11. slope - the slope of the peak exercise ST segment (2 = upsloping; 1 = flat; 0 = downsloping)
12. ca - number of major vessels (0-3) colored by fluoroscopy
13. thal - 2 = normal; 1 = fixed defect; 3 = reversible defect
14. num - the predicted attribute - diagnosis of heart disease (angiographic disease status)
(Value 0 = < diameter narrowing; Value 1 = > 50% diameter narrowing)

- **Data Describe**

```
1 #Numerical Data
2 Heart_Data.describe()
```

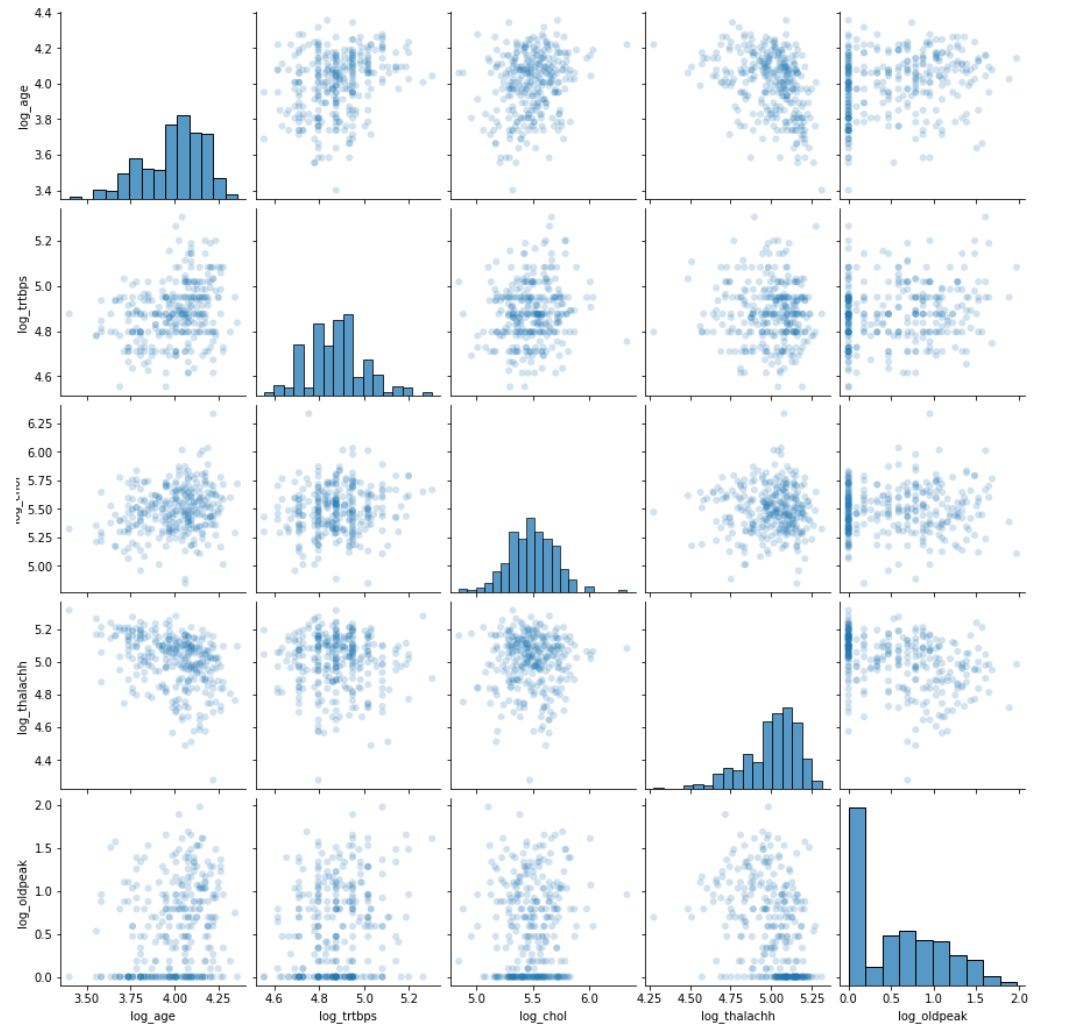
	age	sex	trestbps	chol	fbs	thalachh	exng	oldpeak	caa	output	restecg_hypetrophic	restecg_normal	restecg_abnormal	cp_asymptomatic
count	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000
mean	54.366337	0.683168	131.623762	246.264026	0.148515	149.646865	0.326733	1.039604	0.729373	0.544554	0.485149	0.501650	0.013201	0.471947
std	9.082101	0.466011	17.538143	51.830751	0.356198	22.905161	0.469794	1.161075	1.022606	0.498835	0.500606	0.500824	0.114325	0.500038
min	29.000000	0.000000	94.000000	126.000000	0.000000	71.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	47.500000	0.000000	120.000000	211.000000	0.000000	133.500000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
50%	55.000000	1.000000	130.000000	240.000000	0.000000	153.000000	0.000000	0.800000	0.000000	1.000000	0.000000	1.000000	0.000000	0.000000
75%	61.000000	1.000000	140.000000	274.500000	0.000000	166.000000	1.000000	1.600000	1.000000	1.000000	1.000000	1.000000	0.000000	1.000000
max	77.000000	1.000000	200.000000	564.000000	1.000000	202.000000	1.000000	6.200000	4.000000	1.000000	1.000000	1.000000	1.000000	1.000000

● Numerical Data

Numerical Data Distribution.



The pair plot doesn't show any strong correlation between all numerical data



- **Train / Test Split**

As first step we proceed to split the dataset in Train and Test Split and the using a k-fold cross validation in order to prevent overfitting.

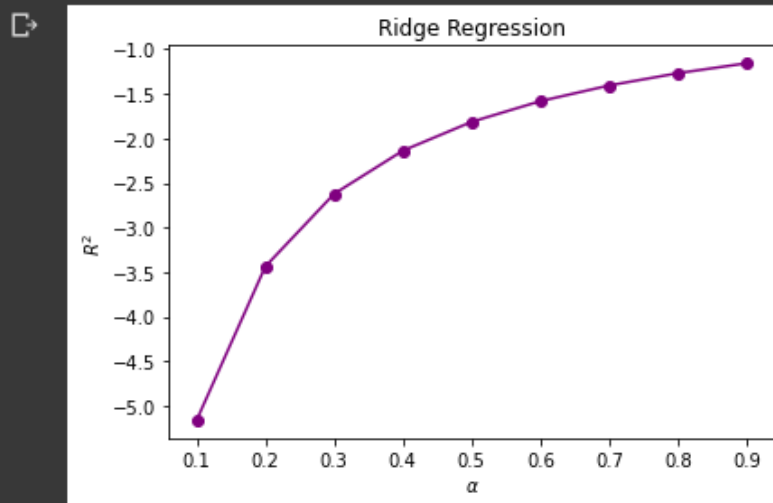
```
10 print(y)
11 #K-Fold
12 kf = KFold(shuffle=True, random_state=72018, n_splits=3)
13
```

- **Finding The Best Model**

Now we proceed to test some regression model to find witch better perform with our data set

Basic Regression

Ridge Regression



```
[ ] 1
    2 best_estimator = Pipeline([
    3 | | | | | | | | | | ("scaler", s),
    4 | | | | | | | | | | ("make_higher_degree", PolynomialFeatures(degree=2)),
    5 | | | | | | | | | | ("ridge_regression", Ridge(alpha=0.03))])
    6
    7 best_estimator.fit(X_train, y_train)
    8 ridge_score = best_estimator.score(X_train, y_train)
```

● Final Thoughts

	linear	lasso	ridge
score	0.489702	0.642254	0.899014

Due to the high number of parameters the ridge regression performs better in estimate . Probably this is more a classification problem so another algorithm can get more precise and significant results