

## 2.0 Notebook 2: Exploratory Data Analysis

Now that we have a good intuitive sense of the data, Next step involves taking a closer look at attributes and data values. In this section, I am getting familiar with the data, which will provide useful knowledge for data pre-processing.

### 2.1 Objectives of Data Exploration

Exploratory data analysis (EDA) is a very important step which takes place after feature engineering and acquiring data and it should be done before any modeling. This is because it is very important for a data scientist to be able to understand the nature of the data without making assumptions. The results of data exploration can be extremely useful in grasping the structure of the data, the distribution of the values, and the presence of extreme values and interrelationships within the data set.

#### The purpose of EDA is:

- to use summary statistics and visualizations to better understand data, \*find clues about the tendencies of the data, its quality and to formulate assumptions and the hypothesis of our analysis
- For data preprocessing to be successful, it is essential to have an overall picture of your data Basic statistical descriptions can be used to identify properties of the data and highlight which data values should be treated as noise or outliers.\*\*

Next step is to explore the data. There are two approached used to examine the data using:

1. **Descriptive statistics** is the process of condensing key characteristics of the data set into simple numeric metrics. Some of the common metrics used are mean, standard deviation, and correlation.
2. **Visualization** is the process of projecting the data, or parts of it, into Cartesian space or into abstract images. In the data mining process, data exploration is leveraged in many different steps including preprocessing, modeling, and interpretation of results.

## 2.2 Descriptive statistics

Summary statistics are measurements meant to describe data. In the field of descriptive statistics, there are many [summary measurements \(http://www.saedsayad.com/numerical\\_variables.htm\)](http://www.saedsayad.com/numerical_variables.htm)

In [1]:

```
%matplotlib inline
import matplotlib.pyplot as plt

#Load libraries for data processing
import pandas as pd #data processing, CSV file I/O (e.g. pd.read_csv)
import numpy as np
from scipy.stats import norm
import seaborn as sns # visualization

plt.rcParams['figure.figsize'] = (15,8)
plt.rcParams['axes.titlesize'] = 'large'
```

In [2]:

```
data = pd.read_csv('data/clean-data.csv', index_col=False)
data.drop('Unnamed: 0',axis=1, inplace=True)
#data.head(2)
```

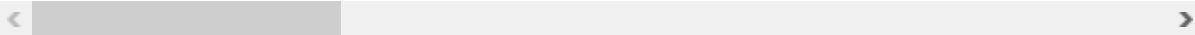
In [3]:

```
#basic descriptive statistics
data.describe()
```

Out[3]:

	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness
count	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000
mean	14.127292	19.289649	91.969033	654.889104	0.096360	0.014064
std	3.524049	4.301036	24.298981	351.914129	0.014064	0.052630
min	6.981000	9.710000	43.790000	143.500000	0.052630	0.086370
25%	11.700000	16.170000	75.170000	420.300000	0.086370	0.095870
50%	13.370000	18.840000	86.240000	551.100000	0.095870	0.105300
75%	15.780000	21.800000	104.100000	782.700000	0.105300	0.163400
max	28.110000	39.280000	188.500000	2501.000000	0.163400	

8 rows × 30 columns



In [4]:

```
data.skew()
```

Out[4]:

```
radius_mean      0.942380
texture_mean     0.650450
perimeter_mean   0.990650
area_mean        1.645732
smoothness_mean  0.456324
compactness_mean 1.190123
concavity_mean   1.401180
concave points_mean 1.171180
symmetry_mean    0.725609
fractal_dimension_mean 1.304489
radius_se        3.088612
texture_se       1.646444
perimeter_se     3.443615
area_se          5.447186
smoothness_se    2.314450
compactness_se   1.902221
concavity_se     5.110463
concave points_se 1.444678
symmetry_se      2.195133
fractal_dimension_se 3.923969
radius_worst     1.103115
texture_worst    0.498321
perimeter_worst  1.128164
area_worst       1.859373
smoothness_worst 0.415426
compactness_worst 1.473555
concavity_worst  1.150237
concave points_worst 0.492616
symmetry_worst   1.433928
fractal_dimension_worst 1.662579
dtype: float64
```

The skew result show a positive (right) or negative (left) skew. Values closer to zero show less skew.

Check binary encoding from NB1 to confirm the conversion of the diagnosis categorical data into numeric, where

- Malignant = 1 (indicates prescence of cancer cells)
- Benign = 0 (indicates abscence)

### Observation

*357 observations indicating the absence of cancer cells and 212 show absence of cancer cell*

Lets confirm this, by plotting the histogram

## 2.3 Data Visualizations

One of the main goals of visualizing the data here is to observe which features are most helpful in predicting malignant or benign cancer. The other is to see general trends that may aid us in model selection and hyper parameter selection.

Apply 3 techniques that you can use to understand each attribute of your dataset independently.

- Histograms.
- Density Plots.

## 2.3.1 Visualise distribution of data via histograms

Histograms are commonly used to visualize numerical variables. A histogram is similar to a bar graph after the values of the variable are grouped (binned) into a finite number of intervals (bins).

Histograms group data into bins and provide you a count of the number of observations in each bin. From the shape of the bins you can quickly get a feeling for whether an attribute is Gaussian, skewed or even has an exponential distribution. It can also help you see possible outliers.

### Separate columns into smaller dataframes to perform visualization

In [8]:

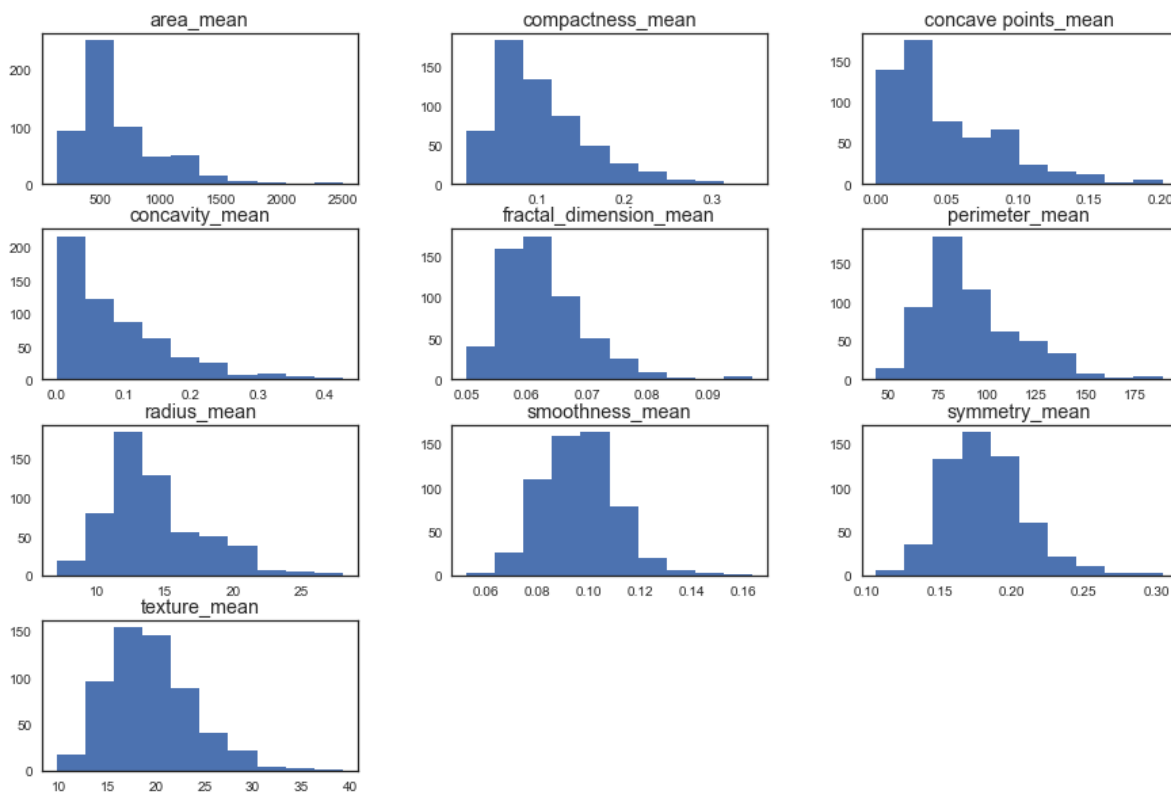
```
#Break up columns into groups, according to their suffix designation
#(_mean, _se,
# and __worst) to perform visualisation plots off.
#Join the 'ID' and 'Diagnosis' back on
data_id_diag=data.loc[:,["id","diagnosis"]]
data_diag=data.loc[:,["diagnosis"]]

#For a merge + slice:
data_mean=data.ix[:,1:11]
```

## Histogram

In [9]:

```
#Plot histograms of CUT1 variables
hist_mean=data_mean.hist(bins=10, figsize=(15, 10),grid=False,)
```



## Observation

We can see that perhaps the attributes **concavity**, and **concavity\_point** may have an exponential distribution ( ). We can also see that perhaps the texture and smooth and symmetry attributes may have a Gaussian or nearly Gaussian distribution. This is interesting because many machine learning techniques assume a Gaussian univariate distribution on the input variables.

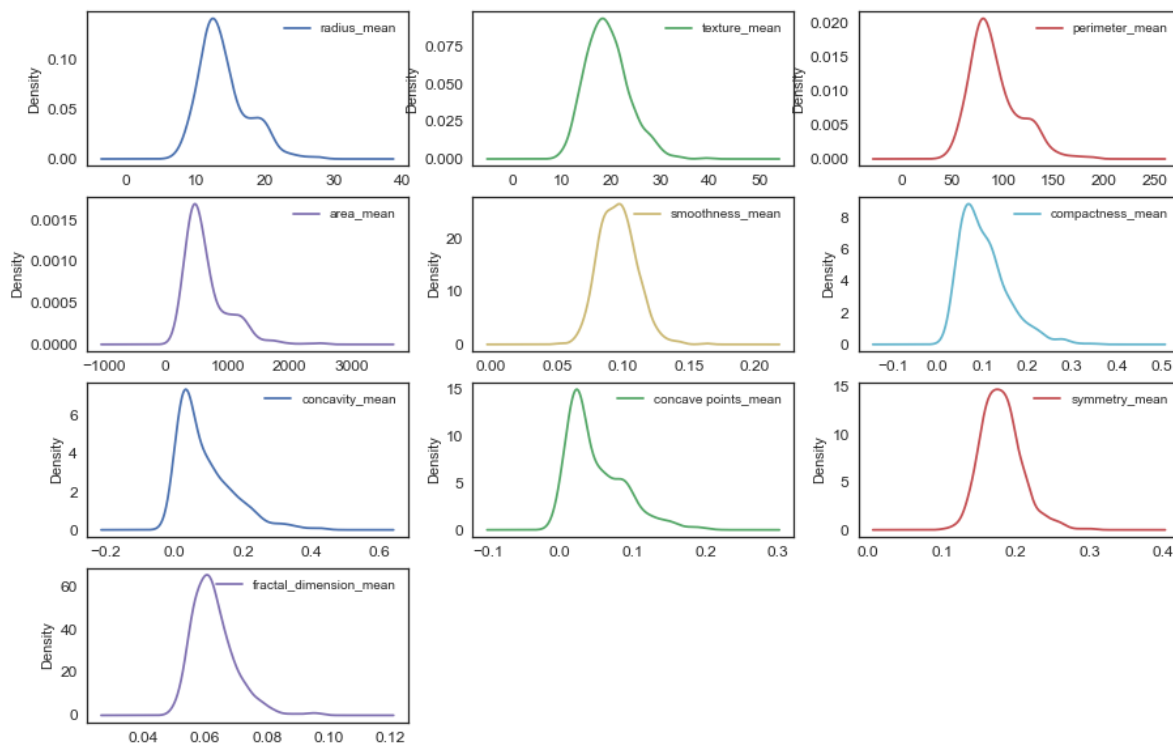
## 2.3.2 Visualize distribution of data via density plots

### Density plots "\_mean" suffix designation

In [12]:

```
#Density Plots
```

```
plt = data_mean.plot(kind= 'density', subplots=True, layout=(4,3), sharex=False,
                      sharey=False, fontsize=12, figsize=(15,10))
```



## Observation

We can see that perhaps the attributes perimeter, radius, area, concavity, compactness may have an exponential distribution ( ). We can also see that perhaps the texture and smooth and symmetry attributes may have a Gaussian or nearly Gaussian distribution. This is interesting because many machine learning techniques assume a Gaussian univariate distribution on the input variables.

## Correlation matrix

In [18]:

```
# plot correlation matrix
import pandas as pd
import numpy as np
import seaborn as sns
from matplotlib import pyplot as plt

plt.style.use('fivethirtyeight')
sns.set_style("white")

data = pd.read_csv('data/clean-data.csv', index_col=False)
data.drop('Unnamed: 0',axis=1, inplace=True)
# Compute the correlation matrix
corr = data.mean().corr()

# Generate a mask for the upper triangle
mask = np.zeros_like(corr, dtype=np.bool)
mask[np.triu_indices_from(mask)] = True

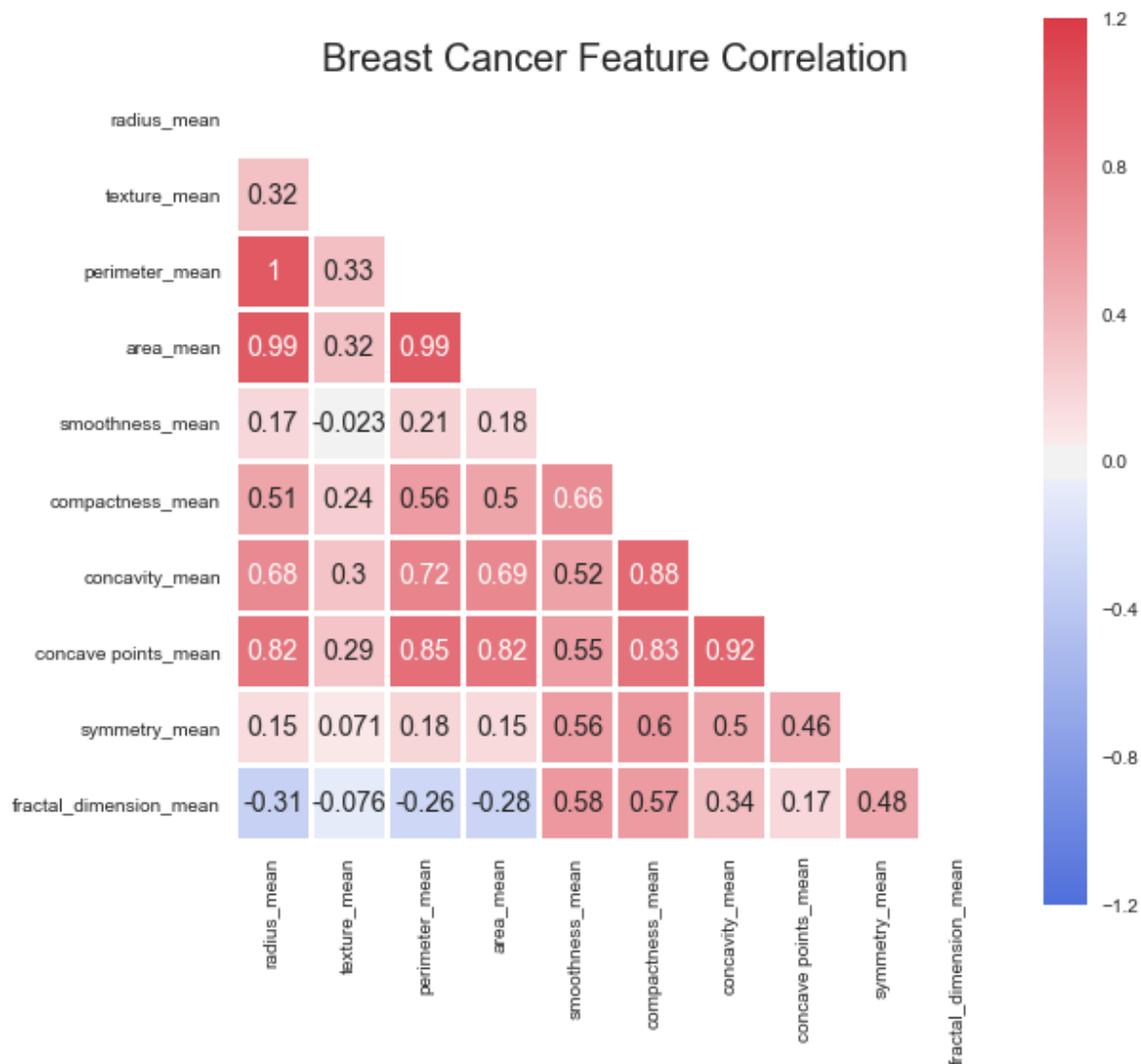
# Set up the matplotlib figure
data, ax = plt.subplots(figsize=(8, 8))
plt.title('Breast Cancer Feature Correlation')

# Generate a custom diverging colormap
cmap = sns.diverging_palette(260, 10, as_cmap=True)

# Draw the heatmap with the mask and correct aspect ratio
sns.heatmap(corr, vmax=1.2, square='square', cmap=cmap, mask=mask,
            ax=ax,annot=True, fmt='.2g',linewidths=2)
```

Out[18]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x114932630>



## Observation:

We can see strong positive relationship exists with mean values paramaters between 1-0.75;.

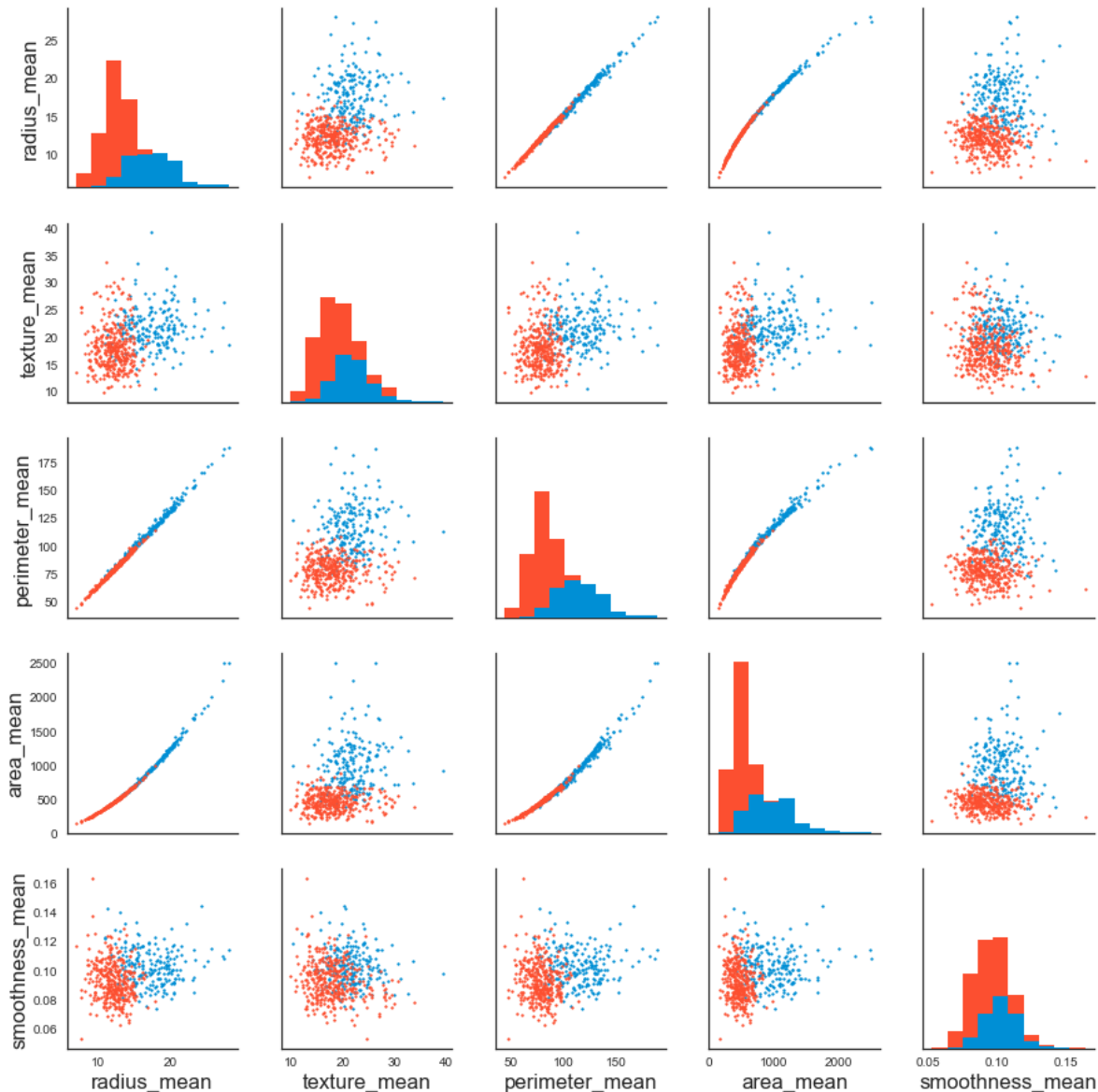
- The mean area of the tissue nucleus has a strong positive correlation with mean values of radius and parameter;
- Some paramters are moderately positive corrlated (r between 0.5-0.75)are concavity and area, concavity and perimeter etc
- Likewise, we see some strong negative correlation between fractal\_dimension with radius, texture, parameter mean values.



In [19]:

```
plt.style.use('fivethirtyeight')
sns.set_style("white")

data = pd.read_csv('data/clean-data.csv', index_col=False)
g = sns.PairGrid(data[[data.columns[1], data.columns[2], data.columns[3],
                      data.columns[4], data.columns[5], data.columns[6]]], hue='diagnosis' )
g = g.map_diag(plt.hist)
g = g.map_offdiag(plt.scatter, s = 3)
```



## Summary

- Mean values of cell radius, perimeter, area, compactness, concavity and concave points can be used in classification of the cancer. Larger values of these parameters tends to show a correlation with malignant tumors.
- mean values of texture, smoothness, symmetry or fractual dimension does not show a particular preference of one diagnosis over the other.
- In any of the histograms there are no noticeable large outliers that warrants further cleanup.

