	<pre>import libraries and dataset import numpy as np import pandas as pd import matplotlib.pyplot as plt from matplotlib import style</pre>
n [2]:	<pre># Importing the dataset, downloaded file is in the same folder csv_path = "cell_samples.csv" cell_samp_df = pd.read_csv(csv_path)</pre>
n [3]:	Run a quick check on the dataset cell_samp_df.shape (699, 11)
n [4]:	<pre>cell_samp_df.info() <class 'pandas.core.frame.dataframe'=""></class></pre>
	Range Index: 699 entries, 0 to 698 Data columns (total 11 columns) Range Index: 699 entries, 0 to 698 Bata columns (total 11 columns) Range Index: 699 entries, 0 to 698 Range Index: 699 entries, 0 to 698 Bata columns (total 11 columns) Range Index: 699 entries, 0 to 698 Bata Columns (total 11 columns) Range Index: 699 entries, 0 to 698 Range Index: 699 entries, 0 to 699 Range
n [5]: ut[5]: _	cell_samp_df.head(5) ID Clump UnifSize UnifShape MargAdh SingEpiSize BareNuc BlandChrom NormNucl Mit Class
	0 1000025 5 1 1 1 2 1 3 1 1 2 1 1002945 5 4 4 5 7 10 3 2 1 2 2 1015425 3 1 1 1 2 2 3 1 1 2 3 1016277 6 8 8 1 3 4 3 7 1 2 4 1017023 4 1 3 2 1 3 1 1 2
	Cell_samp_df.dest=tie() 10
	1
	So, we will treat this as missing data, and we'll replace it using the mean of the column Replace missing values #Replace '?' dengan NaN cell_samp_df["BareNuc"].replace('?', np.NaN, inplace = True)
n [9]:	<pre>#Calculate mean while ignoring NaN value barenuc_mean = int(np.nanmean(cell_samp_df["BareNuc"].astype('float'))) #Convert float (mean) to integer #Replacing the NaN value with the mean cell_samp_df["BareNuc"].replace(np.NaN, barenuc_mean, inplace = True) #Make the type of the column to int again cell_samp_df["BareNuc"] = cell_samp_df["BareNuc"].astype('int')</pre>
[10]: t[10]:	<pre>cell_samp_df["BareNuc"].value_counts() 1 402 10 132</pre>
	3 44 2 30 5 30 8 21 4 19 9 9 7 8 6 4 Name: BareNuc, dtype: int64 As we can see, value_counts for '3' has gone up by 16, meaning we have successfully replaced '?' value to the mean which is '3'
(Omit unnecessary column Column 'ID' is not useful for what we want to do, column 'ID' is just used to identify each patients, but it is not useful for predicting the cancer class (benign or malignant) So, we must omit the column 'ID' to improve the model
[11]:	<pre>cell_samp_df.drop(columns = 'ID', inplace=True) cell_samp_df.head(5)</pre>
	Clump UnifSize UnifShape MargAdh SingEpiSize BareNuc BlandChrom NormNucl Mit Class 1
[13]:	<pre>Splitting dataset #Extracting independent variables: X = np.array(cell_samp_df.drop(columns = 'Class')) X</pre>
t[13]:	array([[5, 1, 1,, 3, 1, 1],
	<pre>y = np.array(cell_samp_df["Class"]) y array([2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 4, 2, 4, 4, 2, 2, 4, 4, 4,</pre>
	4, 2, 4, 2, 2, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4,
	2, 2, 2, 2, 4, 4, 2, 2, 2, 4, 4, 2, 2, 2, 2, 2, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4,
	4, 2, 4, 4, 2, 2, 4, 2, 4, 2, 4, 2, 2, 2, 2, 2, 4, 4, 2, 2, 4, 4, 4, 4, 4, 2, 2, 4, 4, 2, 2, 4, 4, 2, 2, 4, 4, 2, 2, 4, 4, 2, 2, 4, 4, 2, 2, 4, 4, 2, 2, 4, 4, 2, 4, 4, 2, 4, 4, 2, 2, 4, 4, 2, 2, 4, 4, 2, 2, 4, 4, 2, 2, 4, 4, 2, 2, 4, 4, 2, 2, 4, 4, 2, 2, 2, 4, 4, 2, 2, 2, 2, 4, 4, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
[16]: t[16]:	4, 2, 4, 4, 4, 2, 4, 2, 4, 4, 2, 2, 2, 2, 2, 4, 4, 2, 2, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4,
[16]: t[16]:	4, 2, 4, 4, 4, 2, 4, 2, 4, 2, 4, 2, 2, 2, 2, 2, 4, 4, 2, 4, 4, 4, 4, 4, 4, 4, 4, 4, 2, 2, 4, 4, 2, 4, 2, 4, 2, 2, 4, 4, 2, 2, 4, 4, 2, 2, 4, 4, 4, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
[16]: t[16]: t[17]:	4, 2, 4, 4, 4, 2, 4, 2, 4, 2, 4, 2, 2, 2, 2, 2, 4, 2, 2, 4, 4, 4, 4, 4, 4, 4, 2, 2, 4, 2, 2, 4, 2, 2, 4, 4, 2, 4, 2, 2, 4, 4, 2, 2, 4, 4, 2, 2, 4, 4, 2, 2, 4, 4, 2, 2, 4, 4, 2, 4, 2, 2, 4, 2, 2, 2, 4, 2, 2, 2, 4, 2, 2, 2, 4, 2, 2, 2, 4, 2, 2, 2, 4, 2, 2, 2, 4, 2, 2, 2, 4, 2, 2, 2, 4, 2, 2, 2, 4, 2, 2, 2, 4, 2, 2, 2, 4, 2, 2, 2, 4, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
[16]: t[16]: t[17]: t[17]:	1
[16]: t[16]: t[17]: t[17]:	The statem super
[16]: t[16]: t[17]: t[17]:	4, 5, 5, 4, 6, 4, 5, 4, 5, 4, 5, 4, 5, 7, 5, 5, 1, 5, 1, 5, 4, 5, 5, 5, 5, 5, 4, 4, 5, 5, 5, 5, 5, 5, 5, 5, 5, 5, 5, 5, 5,