

# SUMESH R - 20104169

## Basic Analysis using NumPy and Pandas

### Import Libraries

```
In [1]: import pandas as pd
```

```
In [2]: import numpy as np
```

### Import Dataset

```
In [3]: data = pd.read_csv("8_BreastCancerPrediction.csv")
```

```
In [4]: display(data)
```

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	symmetry_mean	fractal_dimension_mean	radius_se	texture_se	perimeter_se	area_se	smoothness_se	compactness_se	concavity_se	concave points_se	symmetry_se	fractal_dimension_se	radius_worst	texture_worst	perimeter_worst	area_worst	smoothness_worst	compactness_worst	concavity_worst	concave points_worst	symmetry_worst	fractal_dimension_worst
0	842302	M	17.99	10.38	122.80	1001.0	0.11840	0.08474	0.09600	0.05263	0.08455	0.05763	1.3860	0.00850	0.00130	0.00070	0.00040	0.00020	0.00010	0.00005	0.00002	0.00001	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000		
1	842517	M	20.57	17.77	132.90	1326.0	0.08474	0.09600	0.05263	0.08455	0.05763	0.05763	1.4900	0.00850	0.00130	0.00070	0.00040	0.00020	0.00010	0.00005	0.00002	0.00001	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000		
2	84300903	M	19.69	21.25	130.00	1203.0	0.10960	0.09600	0.05263	0.08455	0.05763	0.05763	1.4340	0.00850	0.00130	0.00070	0.00040	0.00020	0.00010	0.00005	0.00002	0.00001	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000		
3	84348301	M	11.42	20.38	77.58	386.1	0.14250	0.09600	0.05263	0.08455	0.05763	0.05763	1.3860	0.00850	0.00130	0.00070	0.00040	0.00020	0.00010	0.00005	0.00002	0.00001	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000		
4	84358402	M	20.29	14.34	135.10	1297.0	0.10030	0.09600	0.05263	0.08455	0.05763	0.05763	1.4340	0.00850	0.00130	0.00070	0.00040	0.00020	0.00010	0.00005	0.00002	0.00001	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000		
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...			
564	926424	M	21.56	22.39	142.00	1479.0	0.11100	0.09600	0.05263	0.08455	0.05763	0.05763	1.4340	0.00850	0.00130	0.00070	0.00040	0.00020	0.00010	0.00005	0.00002	0.00001	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000		
565	926682	M	20.13	28.25	131.20	1261.0	0.09780	0.09600	0.05263	0.08455	0.05763	0.05763	1.3860	0.00850	0.00130	0.00070	0.00040	0.00020	0.00010	0.00005	0.00002	0.00001	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000		
566	926954	M	16.60	28.08	108.30	858.1	0.08455	0.09600	0.05263	0.08455	0.05763	0.05763	1.3860	0.00850	0.00130	0.00070	0.00040	0.00020	0.00010	0.00005	0.00002	0.00001	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000		
567	927241	M	20.60	29.33	140.10	1265.0	0.11780	0.09600	0.05263	0.08455	0.05763	0.05763	1.4340	0.00850	0.00130	0.00070	0.00040	0.00020	0.00010	0.00005	0.00002	0.00001	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000		
568	92751	B	7.76	24.54	47.92	181.0	0.05263	0.09600	0.05263	0.08455	0.05763	0.05763	1.3860	0.00850	0.00130	0.00070	0.00040	0.00020	0.00010	0.00005	0.00002	0.00001	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000		

569 rows × 33 columns

### To display top 10 rows

```
In [5]: data.head(10)
```

Out[5]:

	<b>id</b>	<b>diagnosis</b>	<b>radius_mean</b>	<b>texture_mean</b>	<b>perimeter_mean</b>	<b>area_mean</b>	<b>smoothness_mean</b>	<b>com</b>
<b>0</b>	842302	M	17.99	10.38	122.80	1001.0	0.11840	
<b>1</b>	842517	M	20.57	17.77	132.90	1326.0	0.08474	
<b>2</b>	84300903	M	19.69	21.25	130.00	1203.0	0.10960	
<b>3</b>	84348301	M	11.42	20.38	77.58	386.1	0.14250	
<b>4</b>	84358402	M	20.29	14.34	135.10	1297.0	0.10030	
<b>5</b>	843786	M	12.45	15.70	82.57	477.1	0.12780	
<b>6</b>	844359	M	18.25	19.98	119.60	1040.0	0.09463	
<b>7</b>	84458202	M	13.71	20.83	90.20	577.9	0.11890	
<b>8</b>	844981	M	13.00	21.82	87.50	519.8	0.12730	
<b>9</b>	84501001	M	12.46	24.04	83.97	475.9	0.11860	

10 rows × 33 columns

## to display last 5 rows

In [6]:

data.tail()

Out[6]:

	<b>id</b>	<b>diagnosis</b>	<b>radius_mean</b>	<b>texture_mean</b>	<b>perimeter_mean</b>	<b>area_mean</b>	<b>smoothness_mean</b>	<b>co</b>
<b>564</b>	926424	M	21.56	22.39	142.00	1479.0	0.11100	
<b>565</b>	926682	M	20.13	28.25	131.20	1261.0	0.09780	
<b>566</b>	926954	M	16.60	28.08	108.30	858.1	0.08455	
<b>567</b>	927241	M	20.60	29.33	140.10	1265.0	0.11780	
<b>568</b>	92751	B	7.76	24.54	47.92	181.0	0.05263	

5 rows × 33 columns

## statistical summary

In [7]:

data.describe()

Out[7]:

	<b>id</b>	<b>radius_mean</b>	<b>texture_mean</b>	<b>perimeter_mean</b>	<b>area_mean</b>	<b>smoothness_mean</b>	<b>com</b>
<b>count</b>	5.690000e+02	569.000000	569.000000	569.000000	569.000000	569.000000	
<b>mean</b>	3.037183e+07	14.127292	19.289649	91.969033	654.889104	0.096360	

	<b>id</b>	<b>radius_mean</b>	<b>texture_mean</b>	<b>perimeter_mean</b>	<b>area_mean</b>	<b>smoothness_mean</b>	<b>compactness_mean</b>
<b>std</b>	1.250206e+08	3.524049	4.301036	24.298981	351.914129	0.014064	
<b>min</b>	8.670000e+03	6.981000	9.710000	43.790000	143.500000	0.052630	
<b>25%</b>	8.692180e+05	11.700000	16.170000	75.170000	420.300000	0.086370	
<b>50%</b>	9.060240e+05	13.370000	18.840000	86.240000	551.100000	0.095870	
<b>75%</b>	8.813129e+06	15.780000	21.800000	104.100000	782.700000	0.105300	
<b>max</b>	9.113205e+08	28.110000	39.280000	188.500000	2501.000000	0.163400	

8 rows × 32 columns

## To print number of elements

In [8]: `data.size`

Out[8]: 18777

## to print number of row and cols

In [9]: `data.shape`

Out[9]: (569, 33)

## to find missing values

In [10]: `data.isna()`

	<b>id</b>	<b>diagnosis</b>	<b>radius_mean</b>	<b>texture_mean</b>	<b>perimeter_mean</b>	<b>area_mean</b>	<b>smoothness_mean</b>	<b>compactness_mean</b>
<b>0</b>	False	False	False	False	False	False	False	False
<b>1</b>	False	False	False	False	False	False	False	False
<b>2</b>	False	False	False	False	False	False	False	False
<b>3</b>	False	False	False	False	False	False	False	False
<b>4</b>	False	False	False	False	False	False	False	False
...	...	...	...	...	...	...	...	...
<b>564</b>	False	False	False	False	False	False	False	False
<b>565</b>	False	False	False	False	False	False	False	False
<b>566</b>	False	False	False	False	False	False	False	False
<b>567</b>	False	False	False	False	False	False	False	False

	<b>id</b>	<b>diagnosis</b>	<b>radius_mean</b>	<b>texture_mean</b>	<b>perimeter_mean</b>	<b>area_mean</b>	<b>smoothness_mean</b>	<b>comp</b>
<b>568</b>	False	False	False	False	False	False	False	False

569 rows × 33 columns

## fill null values with a constant

In [11]: `data.fillna(5)`

	<b>id</b>	<b>diagnosis</b>	<b>radius_mean</b>	<b>texture_mean</b>	<b>perimeter_mean</b>	<b>area_mean</b>	<b>smoothness_mean</b>	<b>comp</b>
<b>0</b>	842302	M	17.99	10.38	122.80	1001.0	0.11840	
<b>1</b>	842517	M	20.57	17.77	132.90	1326.0	0.08474	
<b>2</b>	84300903	M	19.69	21.25	130.00	1203.0	0.10960	
<b>3</b>	84348301	M	11.42	20.38	77.58	386.1	0.14250	
<b>4</b>	84358402	M	20.29	14.34	135.10	1297.0	0.10030	
...	...	...	...	...	...	...	...	...
<b>564</b>	926424	M	21.56	22.39	142.00	1479.0	0.11100	
<b>565</b>	926682	M	20.13	28.25	131.20	1261.0	0.09780	
<b>566</b>	926954	M	16.60	28.08	108.30	858.1	0.08455	
<b>567</b>	927241	M	20.60	29.33	140.10	1265.0	0.11780	
<b>568</b>	92751	B	7.76	24.54	47.92	181.0	0.05263	

569 rows × 33 columns

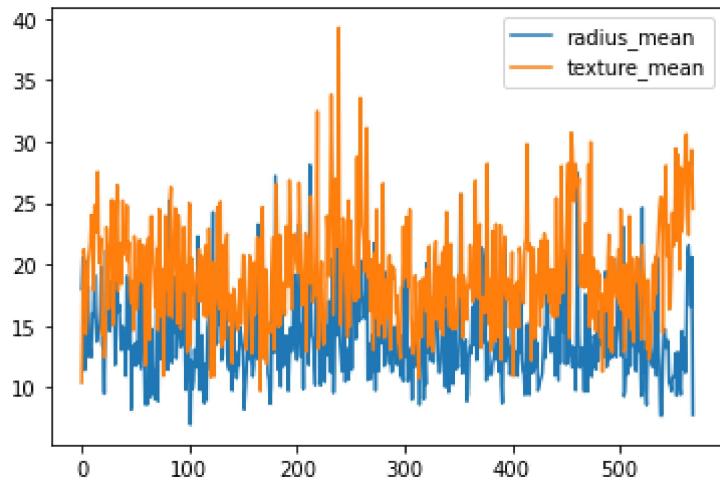
## to select a particular columns

In [12]: `df=pd.DataFrame(data[["radius_mean","texture_mean"]])  
import matplotlib.pyplot as plt`

## line plot

In [13]: `df.plot.line()`

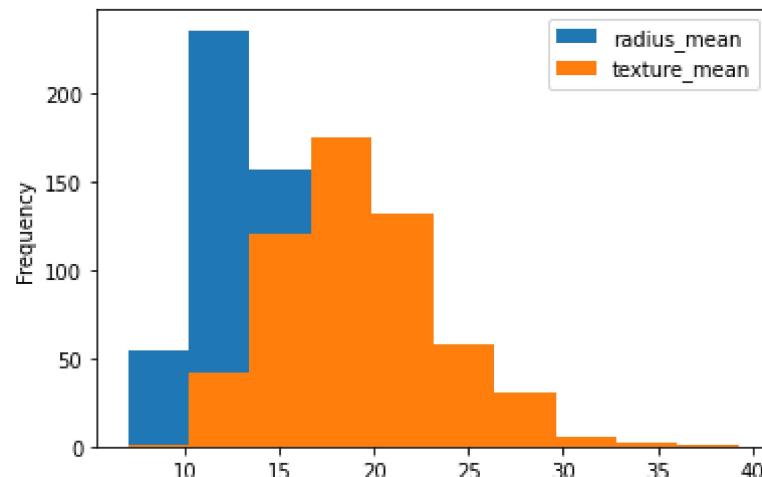
Out[13]: <AxesSubplot:>



## histogram

```
In [14]: df.plot.hist()
```

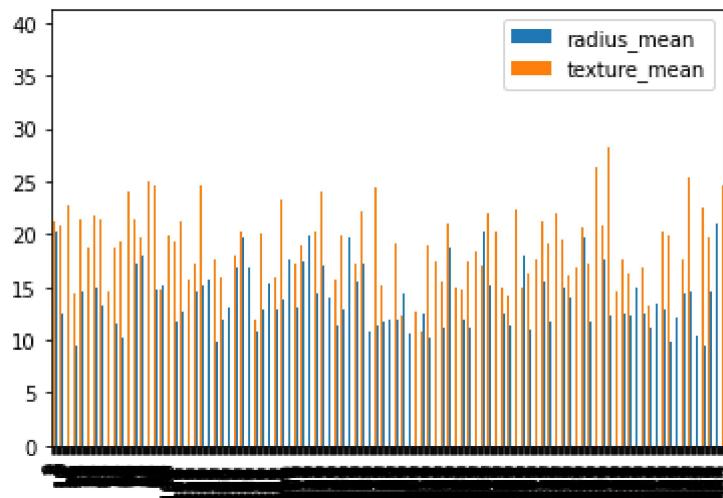
```
Out[14]: <AxesSubplot:ylabel='Frequency'>
```



## bar chart

```
In [15]: df.plot.bar()
```

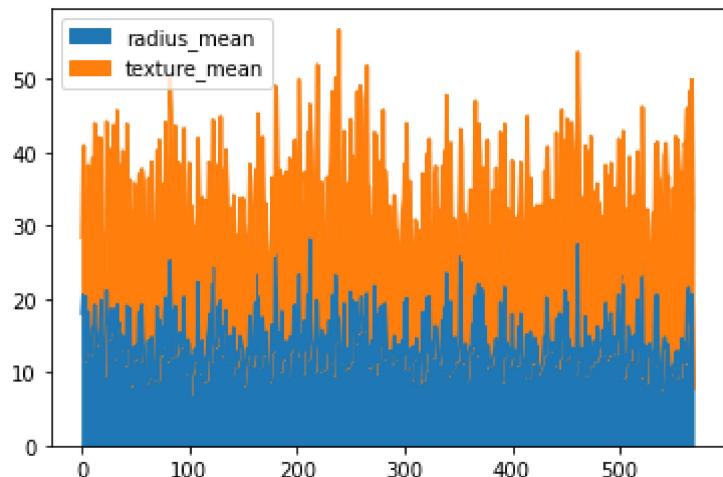
```
Out[15]: <AxesSubplot:>
```



## area plot

```
In [16]: df.plot.area()
```

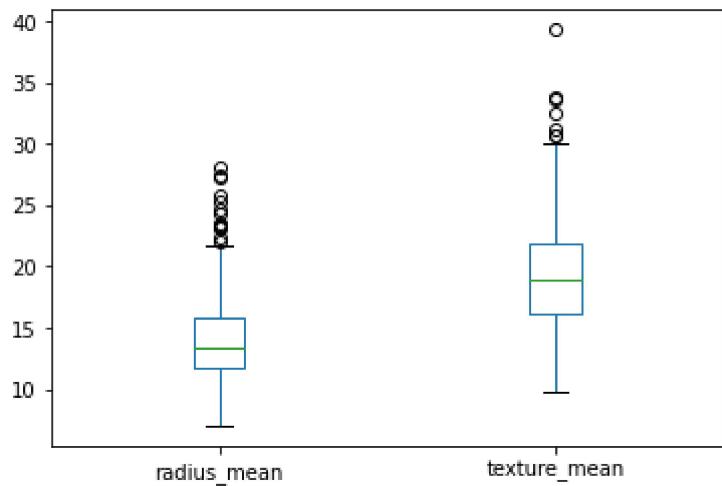
```
Out[16]: <AxesSubplot:>
```



## box plot

```
In [17]: df.plot.box()
```

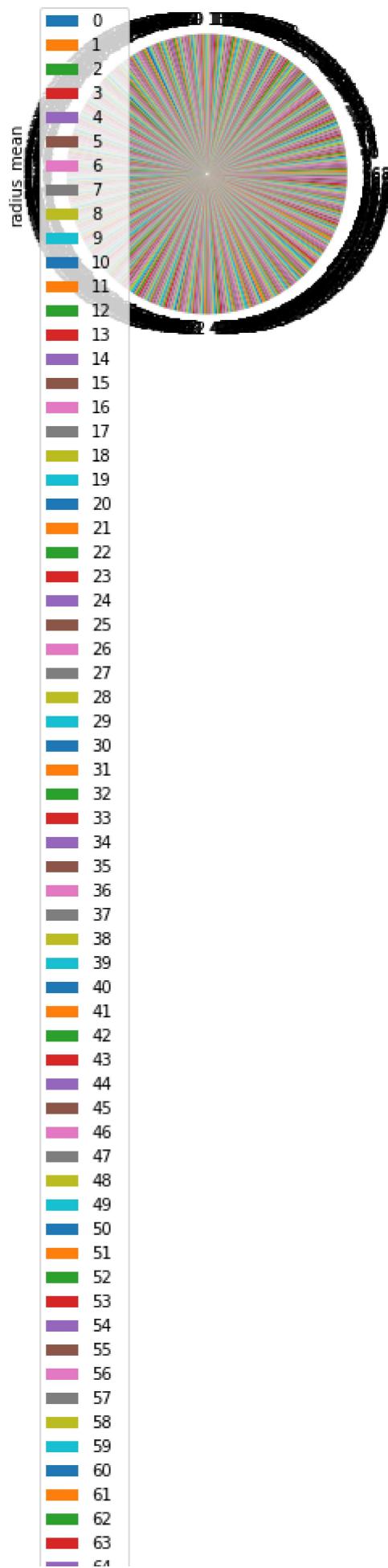
```
Out[17]: <AxesSubplot:>
```

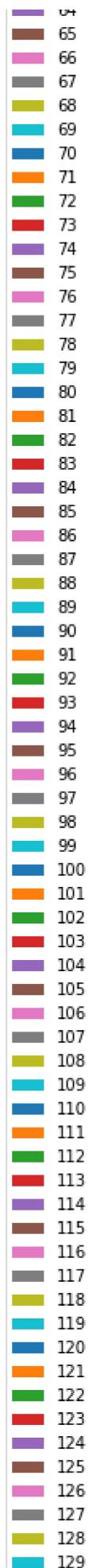


## pie plot

```
In [18]: df.plot.pie(y="radius_mean")
```

```
Out[18]: <AxesSubplot:ylabel='radius_mean'>
```



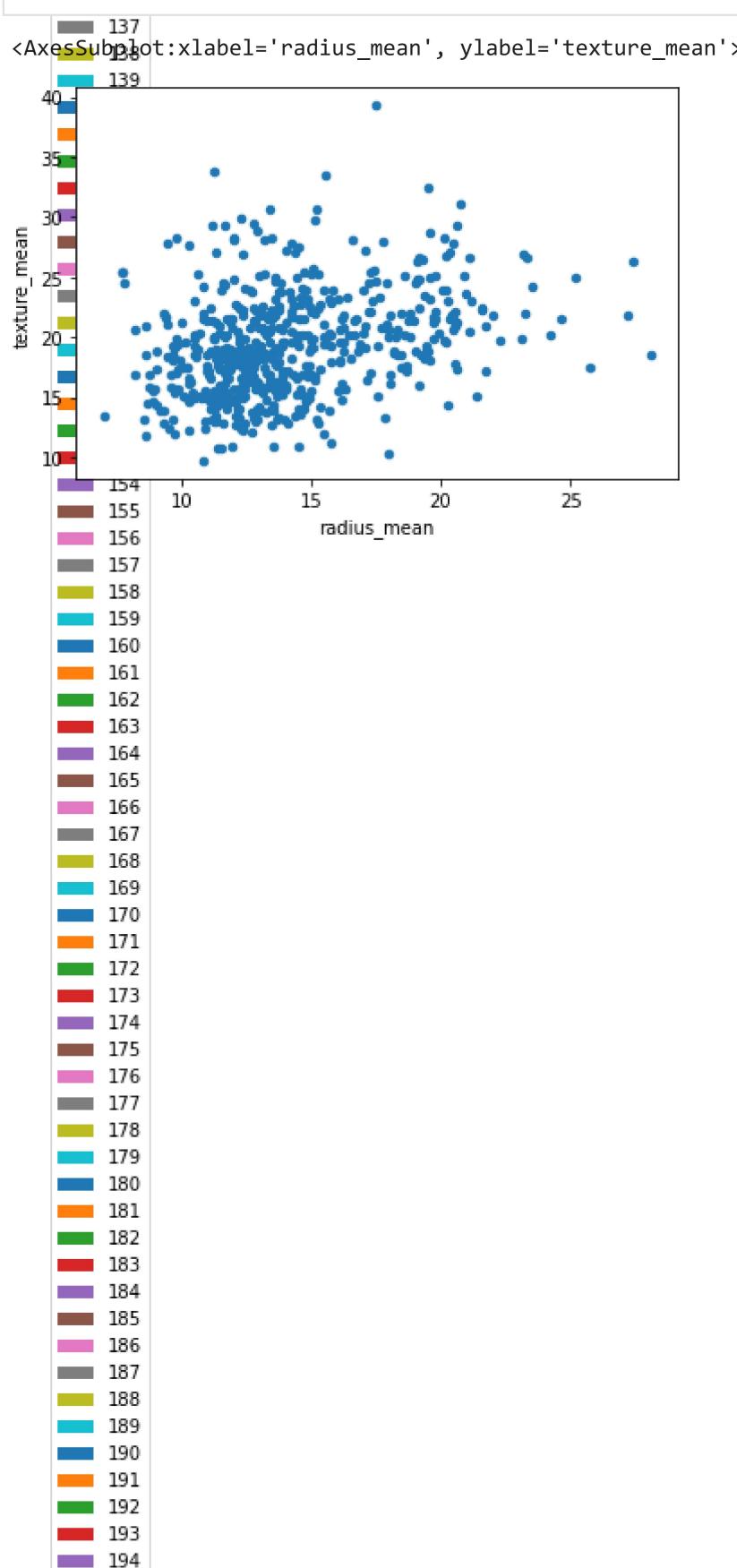


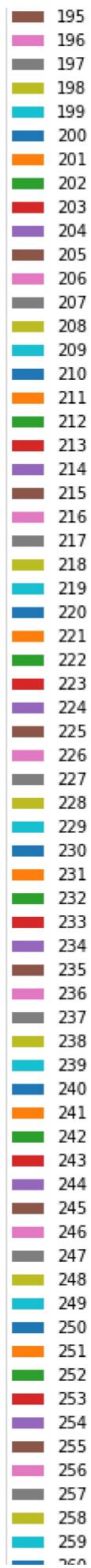
## scatter plot

In [19]:

```
df.plot.scatter(x="radius_mean",y="texture_mean")
```

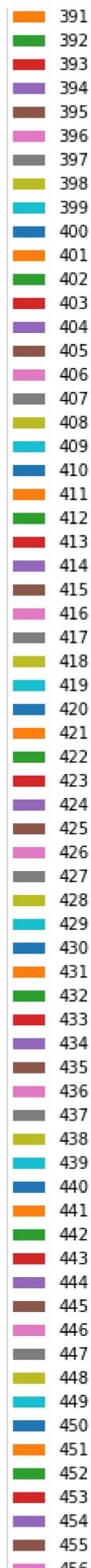
Out[19]:







326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390



450
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521

