# Breast Cancer Wisconsin (Diagnostic)

January 12, 2022

```
[1]: import numpy as np # linear algebra
     import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
     import seaborn as sns  # data visualization library
     import matplotlib.pyplot as plt
     import time
```

```
[2]: data = pd.read_csv('data.csv')
```

```
[3]: data.head()
```

```
[3]:          id diagnosis  radius_mean  texture_mean  perimeter_mean  area_mean  \
     0    842302        M        17.99         10.38          122.80     1001.0
     1    842517        M        20.57         17.77          132.90     1326.0
     2  84300903        M        19.69         21.25          130.00     1203.0
     3  84348301        M        11.42         20.38           77.58      386.1
     4  84358402        M        20.29         14.34          135.10     1297.0

        smoothness_mean  compactness_mean  concavity_mean  concave points_mean  \
     0          0.11840           0.27760          0.3001              0.14710
     1          0.08474           0.07864          0.0869              0.07017
     2          0.10960           0.15990          0.1974              0.12790
     3          0.14250           0.28390          0.2414              0.10520
     4          0.10030           0.13280          0.1980              0.10430

        …  texture_worst  perimeter_worst  area_worst  smoothness_worst  \
     0  …          17.33           184.60      2019.0            0.1622
     1  …          23.41           158.80      1956.0            0.1238
     2  …          25.53           152.50      1709.0            0.1444
     3  …          26.50            98.87       567.7            0.2098
     4  …          16.67           152.20      1575.0            0.1374

        compactness_worst  concavity_worst  concave points_worst  symmetry_worst  \
     0             0.6656           0.7119                0.2654          0.4601
     1             0.1866           0.2416                0.1860          0.2750
     2             0.4245           0.4504                0.2430          0.3613
     3             0.8663           0.6869                0.2575          0.6638
     4             0.2050           0.4000                0.1625          0.2364
```

```
     fractal_dimension_worst   Unnamed: 32
0                    0.11890           NaN
1                    0.08902           NaN
2                    0.08758           NaN
3                    0.17300           NaN
4                    0.07678           NaN

[5 rows x 33 columns]
```

[4]:
```python
col = data.columns
print(col)
```

```
Index(['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', 'Unnamed: 32'],
      dtype='object')
```

[5]:
```python
y = data.diagnosis
drop_cols = ['Unnamed: 32','id','diagnosis']
x = data.drop(drop_cols, axis=1)
x.head()
```

[5]:
```
   radius_mean  texture_mean  perimeter_mean  area_mean  smoothness_mean  \
0        17.99         10.38          122.80     1001.0          0.11840
1        20.57         17.77          132.90     1326.0          0.08474
2        19.69         21.25          130.00     1203.0          0.10960
3        11.42         20.38           77.58      386.1          0.14250
4        20.29         14.34          135.10     1297.0          0.10030

   compactness_mean  concavity_mean  concave points_mean  symmetry_mean  \
0           0.27760          0.3001              0.14710         0.2419
1           0.07864          0.0869              0.07017         0.1812
2           0.15990          0.1974              0.12790         0.2069
3           0.28390          0.2414              0.10520         0.2597
4           0.13280          0.1980              0.10430         0.1809

   fractal_dimension_mean  …  radius_worst  texture_worst  perimeter_worst  \
0                 0.07871  …         25.38          17.33           184.60
1                 0.05667  …         24.99          23.41           158.80
2                 0.05999  …         23.57          25.53           152.50
```

```
3                    0.09744  …        14.91          26.50            98.87
4                    0.05883  …        22.54          16.67           152.20

   area_worst  smoothness_worst  compactness_worst  concavity_worst  \
0     2019.0            0.1622             0.6656           0.7119
1     1956.0            0.1238             0.1866           0.2416
2     1709.0            0.1444             0.4245           0.4504
3      567.7            0.2098             0.8663           0.6869
4     1575.0            0.1374             0.2050           0.4000

   concave points_worst  symmetry_worst  fractal_dimension_worst
0                0.2654          0.4601                  0.11890
1                0.1860          0.2750                  0.08902
2                0.2430          0.3613                  0.08758
3                0.2575          0.6638                  0.17300
4                0.1625          0.2364                  0.07678

[5 rows x 30 columns]
```
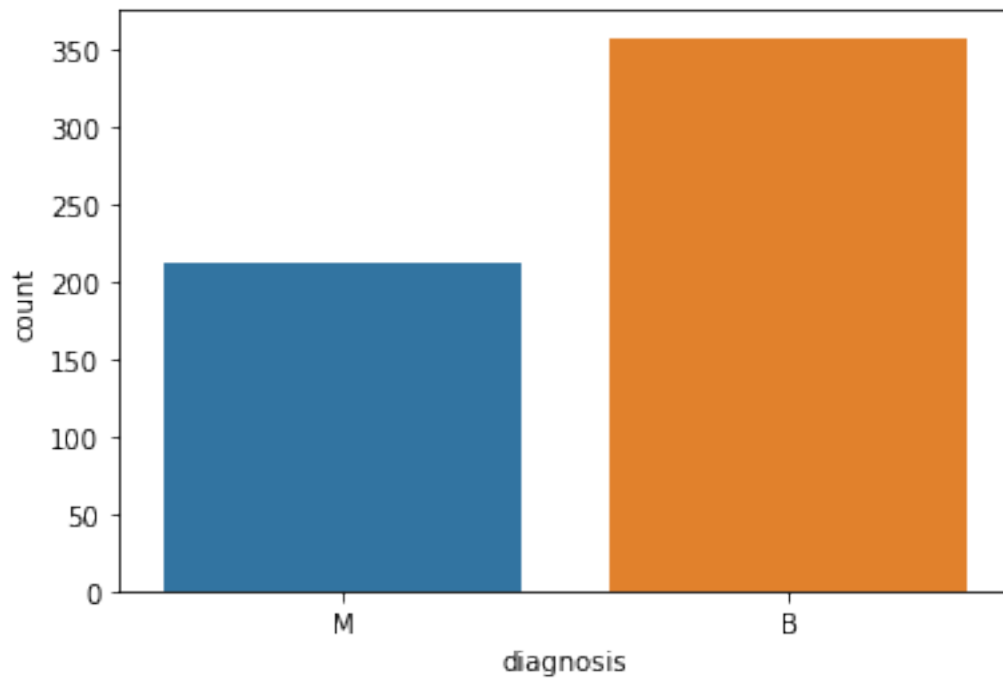
[6]:
```python
ax = sns.countplot(y, label="Count")
B, M = y.value_counts()
print('Number of Benign Tumors', B)
print('Number of Malignant Tumors', M)
```

```
Number of Benign Tumors 357
Number of Malignant Tumors 212
```

```
[7]: x.describe()
```

```
[7]:        radius_mean  texture_mean  perimeter_mean     area_mean  \
     count   569.000000    569.000000      569.000000    569.000000
     mean     14.127292     19.289649       91.969033    654.889104
     std       3.524049      4.301036       24.298981    351.914129
     min       6.981000      9.710000       43.790000    143.500000
     25%      11.700000     16.170000       75.170000    420.300000
     50%      13.370000     18.840000       86.240000    551.100000
     75%      15.780000     21.800000      104.100000    782.700000
     max      28.110000     39.280000      188.500000   2501.000000

            smoothness_mean  compactness_mean  concavity_mean  concave points_mean  \
     count       569.000000        569.000000      569.000000           569.000000
     mean          0.096360          0.104341        0.088799             0.048919
     std           0.014064          0.052813        0.079720             0.038803
     min           0.052630          0.019380        0.000000             0.000000
     25%           0.086370          0.064920        0.029560             0.020310
     50%           0.095870          0.092630        0.061540             0.033500
     75%           0.105300          0.130400        0.130700             0.074000
     max           0.163400          0.345400        0.426800             0.201200

            symmetry_mean  fractal_dimension_mean  …  radius_worst  \
     count     569.000000              569.000000  …    569.000000
     mean        0.181162                0.062798  …     16.269190
     std         0.027414                0.007060  …      4.833242
     min         0.106000                0.049960  …      7.930000
     25%         0.161900                0.057700  …     13.010000
     50%         0.179200                0.061540  …     14.970000
     75%         0.195700                0.066120  …     18.790000
     max         0.304000                0.097440  …     36.040000

            texture_worst  perimeter_worst    area_worst  smoothness_worst  \
     count     569.000000       569.000000    569.000000        569.000000
     mean       25.677223       107.261213    880.583128          0.132369
     std         6.146258        33.602542    569.356993          0.022832
     min        12.020000        50.410000    185.200000          0.071170
     25%        21.080000        84.110000    515.300000          0.116600
     50%        25.410000        97.660000    686.500000          0.131300
     75%        29.720000       125.400000   1084.000000          0.146000
     max        49.540000       251.200000   4254.000000          0.222600

            compactness_worst  concavity_worst  concave points_worst  \
     count         569.000000       569.000000            569.000000
     mean            0.254265         0.272188              0.114606
     std             0.157336         0.208624              0.065732
```
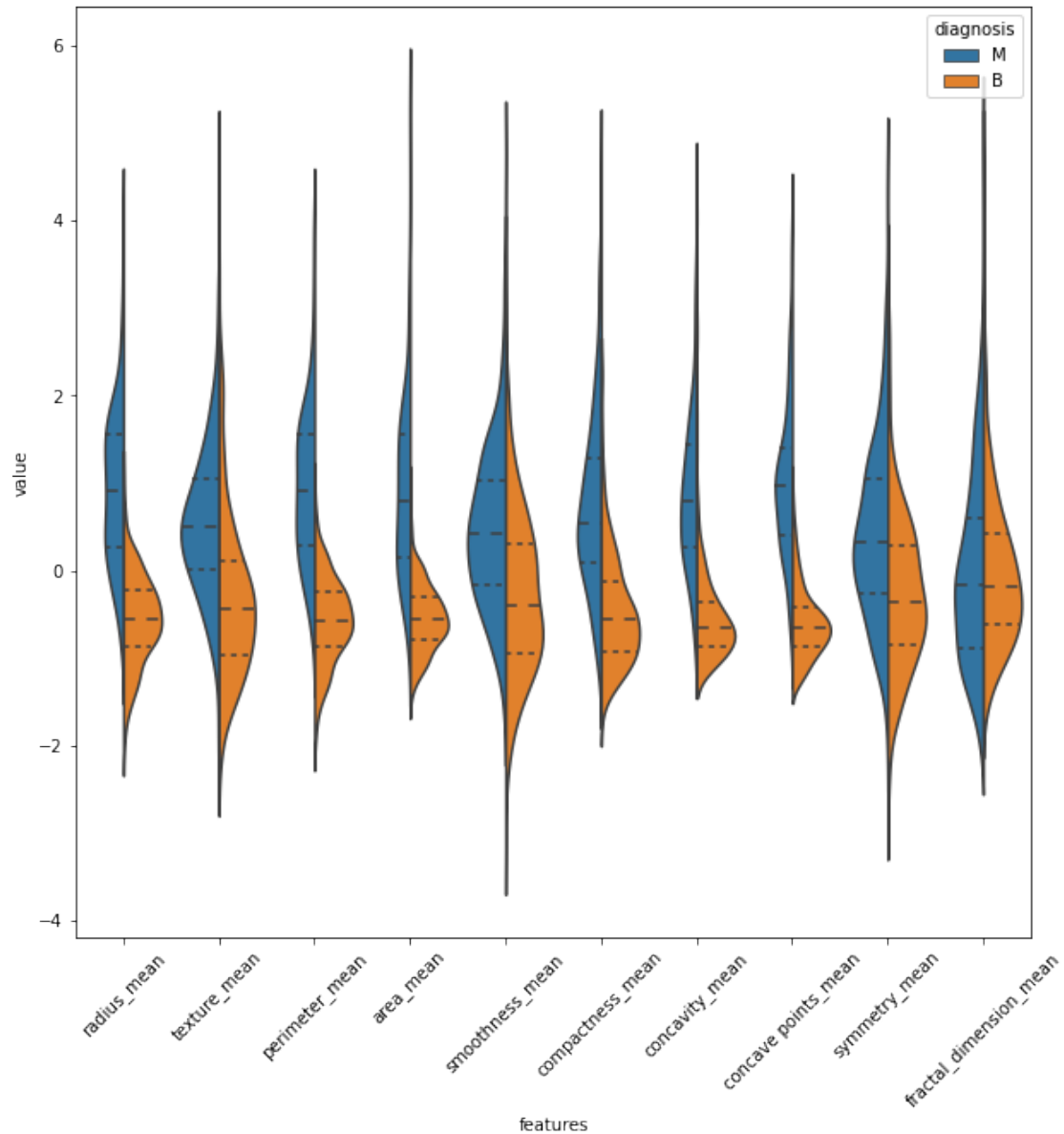
```
min           0.027290      0.000000           0.000000
25%           0.147200      0.114500           0.064930
50%           0.211900      0.226700           0.099930
75%           0.339100      0.382900           0.161400
max           1.058000      1.252000           0.291000


       symmetry_worst  fractal_dimension_worst
count      569.000000               569.000000
mean         0.290076                 0.083946
std          0.061867                 0.018061
min          0.156500                 0.055040
25%          0.250400                 0.071460
50%          0.282200                 0.080040
75%          0.317900                 0.092080
max          0.663800                 0.207500

[8 rows x 30 columns]
```
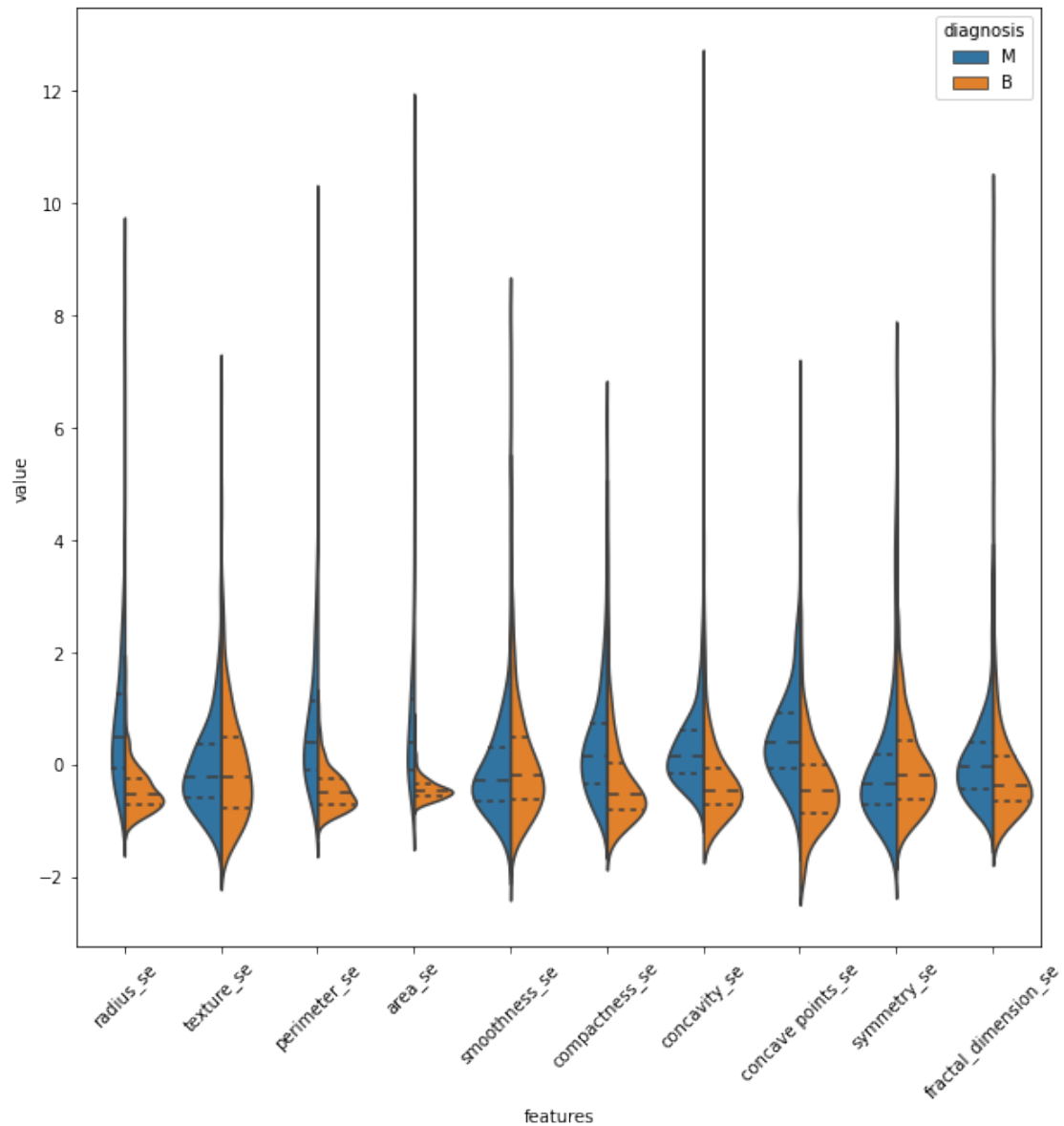
```python
[8]: data = x
     data_std = (data - data.mean()) / data.std()
     data = pd.concat([y, data_std.iloc[:,0:10]], axis=1)
     data = pd.melt(data, id_vars='diagnosis',
                    var_name='features',
                    value_name='value')
     plt.figure(figsize=(10,10))
     sns.violinplot(x='features', y='value', hue='diagnosis', data=data, split=True,
      ↪inner='quart')
     plt.xticks(rotation=45);
```
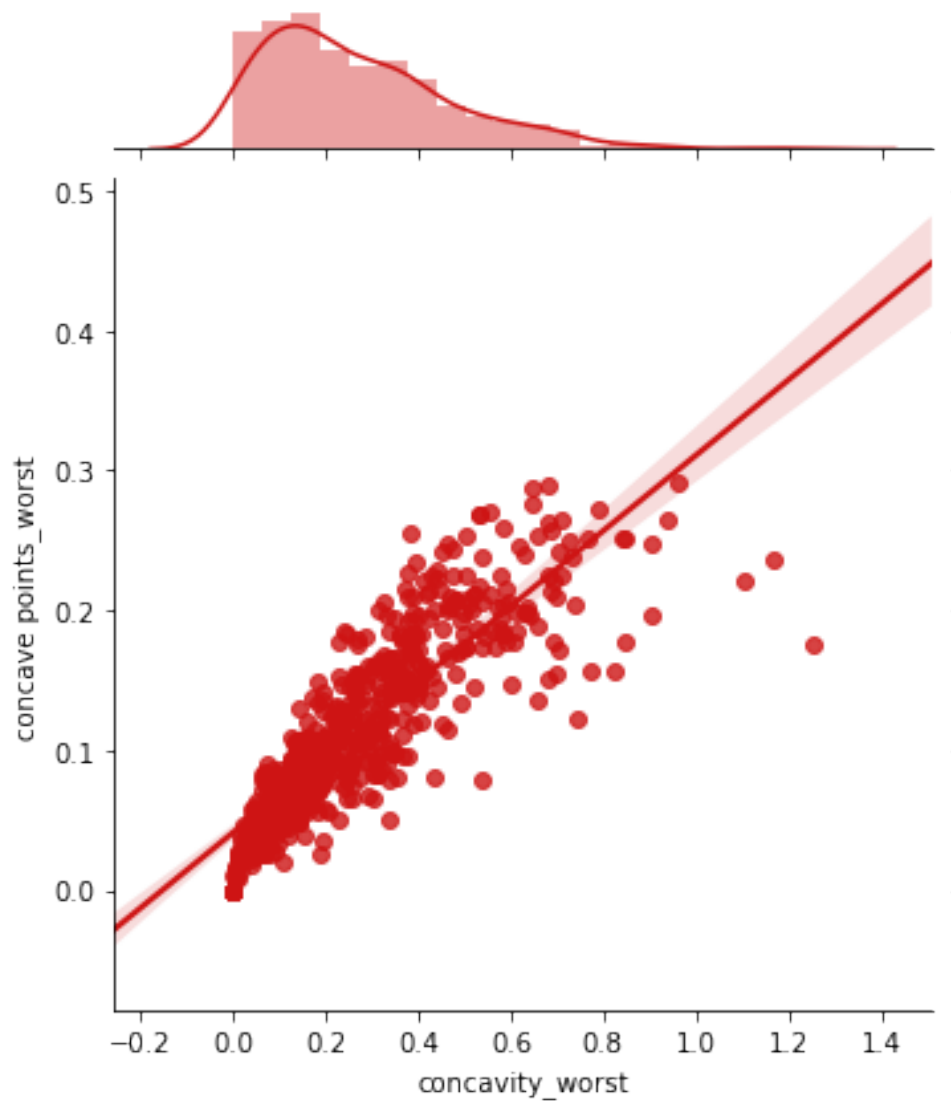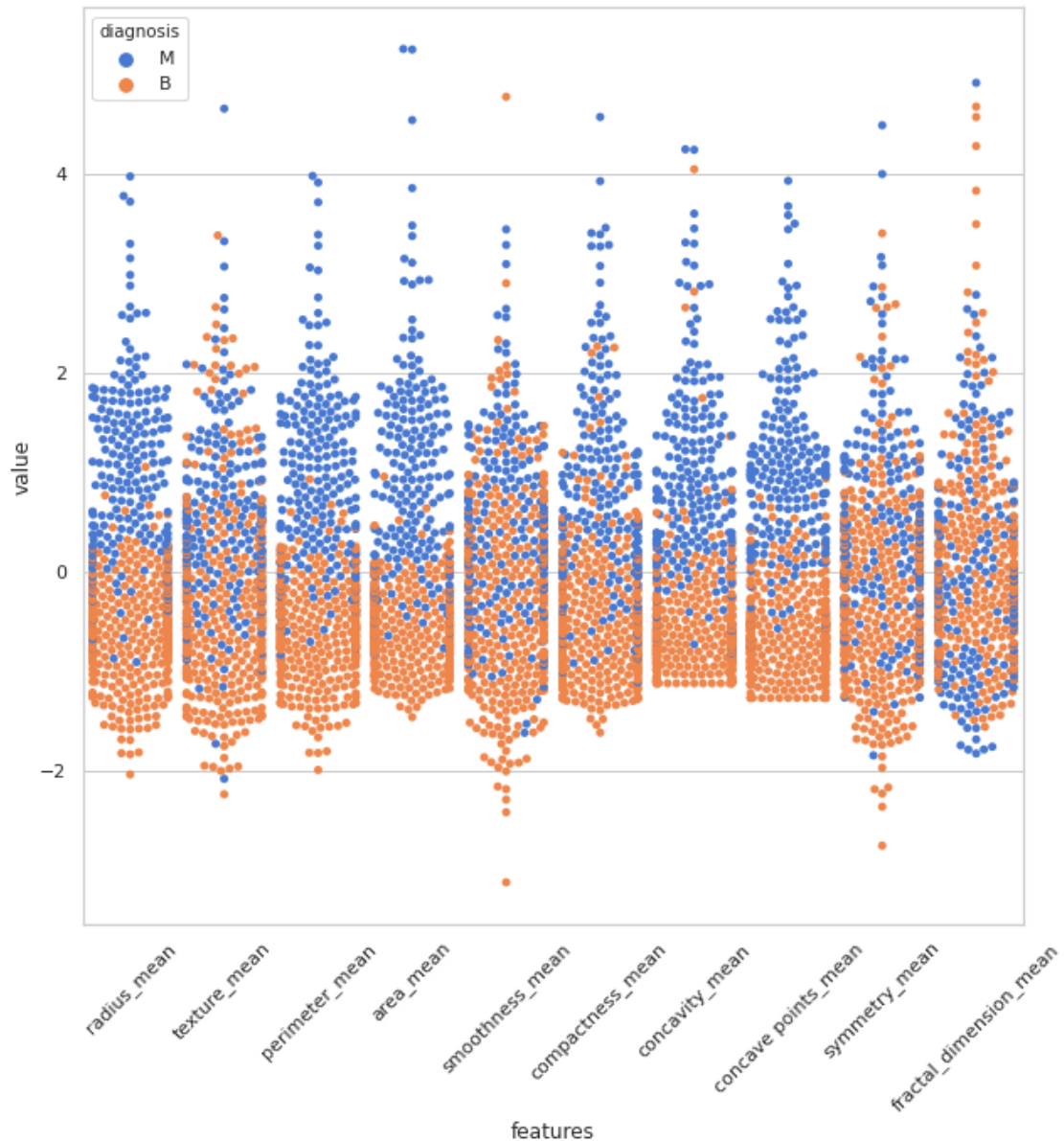
```
[9]: data = pd.concat([y, data_std.iloc[:,10:20]], axis=1)
     data = pd.melt(data, id_vars='diagnosis',
                    var_name='features',
                    value_name='value')
     plt.figure(figsize=(10,10))
     sns.violinplot(x='features', y='value', hue='diagnosis', data=data, split=True,␣
      ↪inner='quart')
     plt.xticks(rotation=45);
```

```
[10]: data = pd.concat([y, data_std.iloc[:,20:30]], axis=1)
      data = pd.melt(data, id_vars='diagnosis',
                     var_name='features',
                     value_name='value')
      plt.figure(figsize=(10,10))
      sns.violinplot(x='features', y='value', hue='diagnosis', data=data, split=True,␣
       ↪inner='quart')
      plt.xticks(rotation=45);
```

```
[11]: sns.boxplot(x='features', y='value', hue='diagnosis', data=data)
      plt.xticks(rotation=45);
```
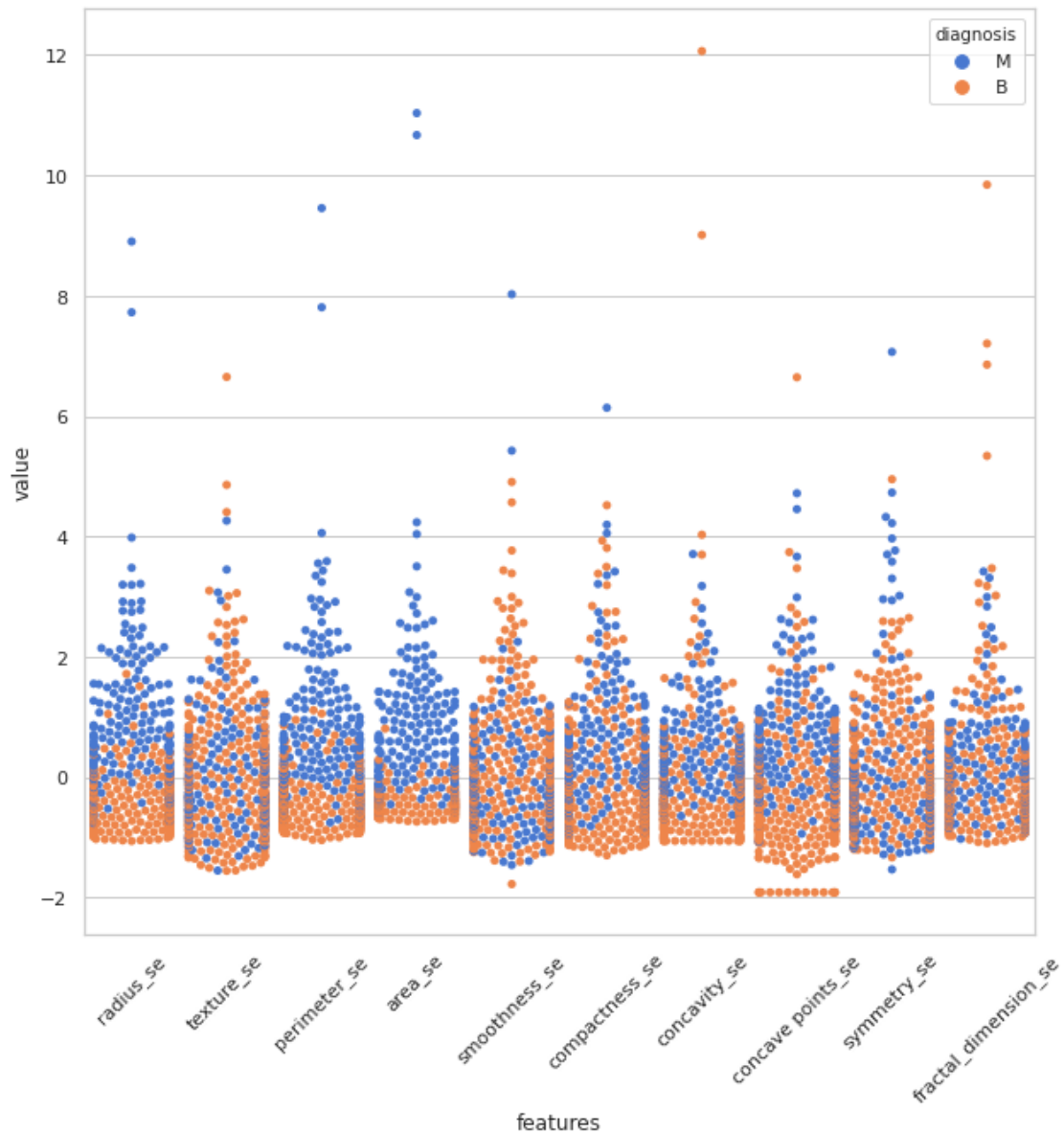
```
[12]: sns.jointplot(x.loc[:, 'concavity_worst'],
                     x.loc[:, 'concave points_worst'],
                     kind='regg',
                     color='#ce1414');
```
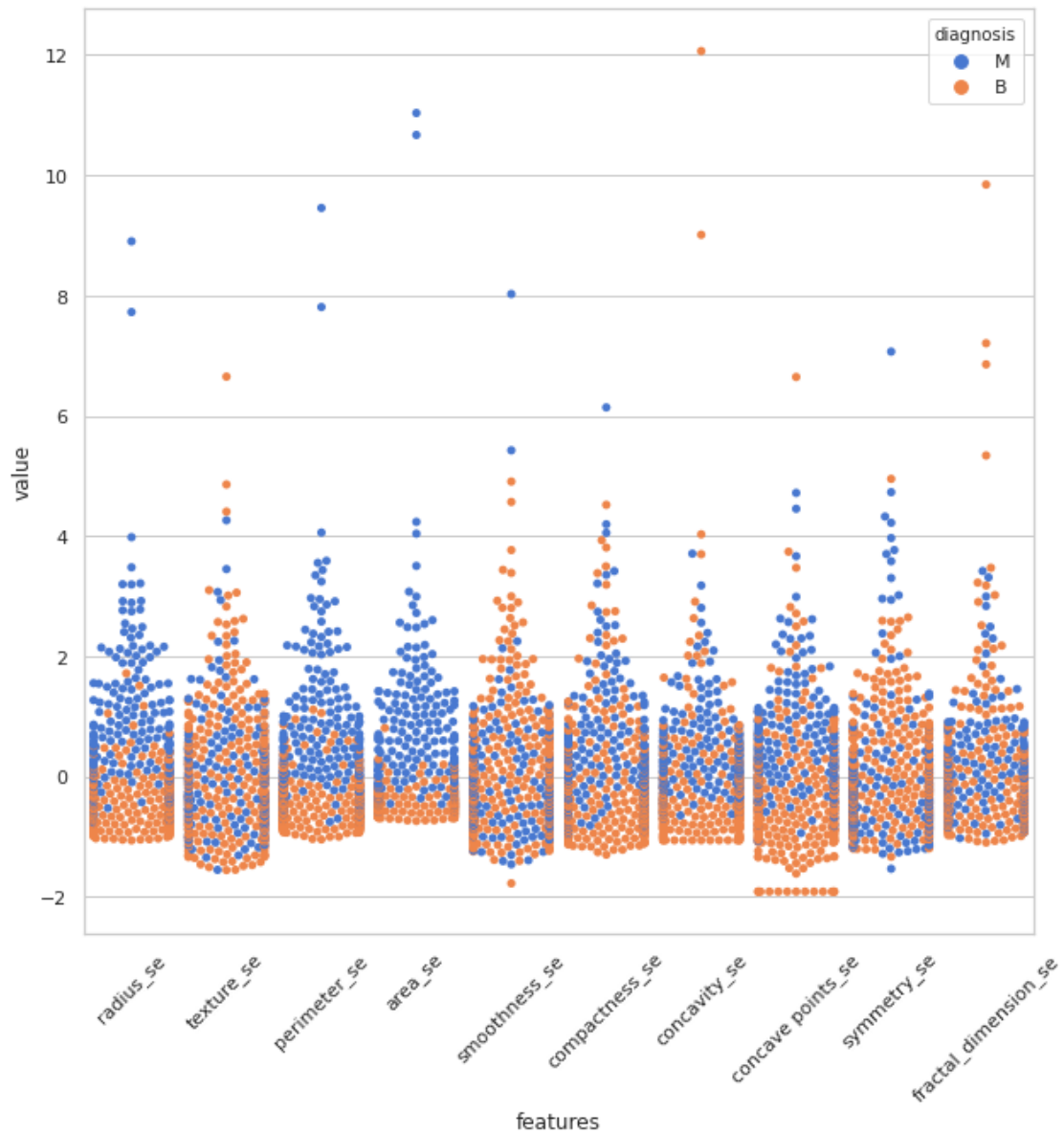
```
sns.set(style='whitegrid', palette='muted')
data = x
data_std = (data - data.mean()) / data.std()
data = pd.concat([y, data_std.iloc[:,0:10]], axis=1)
data = pd.melt(data, id_vars='diagnosis',
              var_name='features',
              value_name='value')
plt.figure(figsize=(10,10))
sns.swarmplot(x='features', y='value', hue='diagnosis', data=data)
plt.xticks(rotation=45);
```

```
[14]: sns.set(style='whitegrid', palette='muted')
      data = x
      data_std = (data - data.mean()) / data.std()
      data = pd.concat([y, data_std.iloc[:,10:20]], axis=1)
      data = pd.melt(data, id_vars='diagnosis',
                  var_name='features',
                  value_name='value')
      plt.figure(figsize=(10,10))
      sns.swarmplot(x='features', y='value', hue='diagnosis', data=data)
      plt.xticks(rotation=45);
```
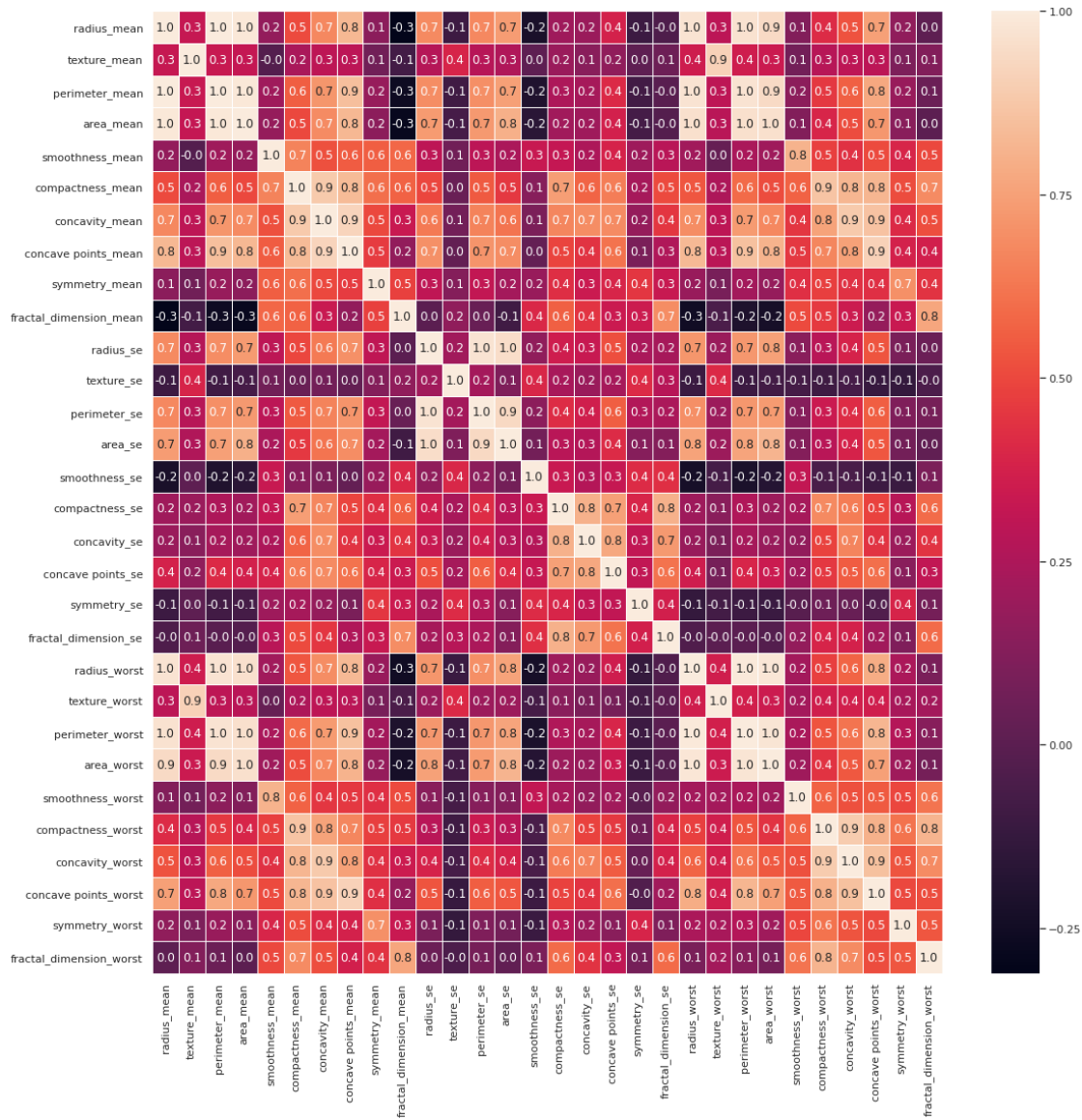
```python
[15]: sns.set(style='whitegrid', palette='muted')
      data = x
      data_std = (data - data.mean()) / data.std()
      data = pd.concat([y, data_std.iloc[:,10:20]], axis=1)
      data = pd.melt(data, id_vars='diagnosis',
                  var_name='features',
                  value_name='value')
      plt.figure(figsize=(10,10))
      sns.swarmplot(x='features', y='value', hue='diagnosis', data=data)
      plt.xticks(rotation=45);
```

```
[16]:  f, ax = plt.subplots(figsize=(18, 18))
       sns.heatmap(x.corr(), annot=True, linewidth=.5, fmt='.1f', ax=ax)
```

[16]: <AxesSubplot:>