**Darshan UNIVERSITY**
योग: कर्मसु कौशलम्

# Data Mining

# Lab - 3

## 1) First, you need to read the titanic dataset from local disk and display first five records

In [2]:
```python
import pandas as pd
```

In [5]:
```python
df=pd.read_csv("titanic.csv")
```

In [6]:
```python
df.head(5)
```

Out[6]:

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | S |
| 1 | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | C |
| 2 | 3 | 1 | 3 | Heikkinen, Miss. Laina | female | 26.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | S |
| 3 | 4 | 1 | 1 | Futrelle, Mrs. Jacques Heath (Lily May Peel) | female | 35.0 | 1 | 0 | 113803 | 53.1000 | C123 | S |
| 4 | 5 | 0 | 3 | Allen, Mr. William Henry | male | 35.0 | 0 | 0 | 373450 | 8.0500 | NaN | S |

In [ ]:

### 2)  Identify Nominal, Ordinal, Binary and Numeric attributes from data sets and display all values.

In [ ]:
```python
numeric=["PassengerId","Age","SibSp","Parch","Fare",]
Nominal=["Name","Ticket","Cabin","Embarked"]
Ordinal=["Pclass"]
Binary=["Survived","Sex"]
```

In [12]: `df["PassengerId"].unique()`

Out[12]:
```
array([  1,   2,   3,   4,   5,   6,   7,   8,   9,  10,  11,  12,  13,
        14,  15,  16,  17,  18,  19,  20,  21,  22,  23,  24,  25,  26,
        27,  28,  29,  30,  31,  32,  33,  34,  35,  36,  37,  38,  39,
        40,  41,  42,  43,  44,  45,  46,  47,  48,  49,  50,  51,  52,
        53,  54,  55,  56,  57,  58,  59,  60,  61,  62,  63,  64,  65,
        66,  67,  68,  69,  70,  71,  72,  73,  74,  75,  76,  77,  78,
        79,  80,  81,  82,  83,  84,  85,  86,  87,  88,  89,  90,  91,
        92,  93,  94,  95,  96,  97,  98,  99, 100, 101, 102, 103, 104,
       105, 106, 107, 108, 109, 110, 111, 112, 113, 114, 115, 116, 117,
       118, 119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129, 130,
       131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143,
       144, 145, 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, 156,
       157, 158, 159, 160, 161, 162, 163, 164, 165, 166, 167, 168, 169,
       170, 171, 172, 173, 174, 175, 176, 177, 178, 179, 180, 181, 182,
       183, 184, 185, 186, 187, 188, 189, 190, 191, 192, 193, 194, 195,
       196, 197, 198, 199, 200, 201, 202, 203, 204, 205, 206, 207, 208,
       209, 210, 211, 212, 213, 214, 215, 216, 217, 218, 219, 220, 221,
       222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232, 233, 234,
       235, 236, 237, 238, 239, 240, 241, 242, 243, 244, 245, 246, 247,
       248, 249, 250, 251, 252, 253, 254, 255, 256, 257, 258, 259, 260,
       261, 262, 263, 264, 265, 266, 267, 268, 269, 270, 271, 272, 273,
       274, 275, 276, 277, 278, 279, 280, 281, 282, 283, 284, 285, 286,
       287, 288, 289, 290, 291, 292, 293, 294, 295, 296, 297, 298, 299,
       300, 301, 302, 303, 304, 305, 306, 307, 308, 309, 310, 311, 312,
       313, 314, 315, 316, 317, 318, 319, 320, 321, 322, 323, 324, 325,
       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,
       391, 392, 393, 394, 395, 396, 397, 398, 399, 400, 401, 402, 403,
       404, 405, 406, 407, 408, 409, 410, 411, 412, 413, 414, 415, 416,
       417, 418, 419, 420, 421, 422, 423, 424, 425, 426, 427, 428, 429,
       430, 431, 432, 433, 434, 435, 436, 437, 438, 439, 440, 441, 442,
       443, 444, 445, 446, 447, 448, 449, 450, 451, 452, 453, 454, 455,
       456, 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, 522, 523, 524, 525, 526, 527, 528, 529, 530, 531, 532, 533,
       534, 535, 536, 537, 538, 539, 540, 541, 542, 543, 544, 545, 546,
       547, 548, 549, 550, 551, 552, 553, 554, 555, 556, 557, 558, 559,
       560, 561, 562, 563, 564, 565, 566, 567, 568, 569, 570, 571, 572,
       573, 574, 575, 576, 577, 578, 579, 580, 581, 582, 583, 584, 585,
       586, 587, 588, 589, 590, 591, 592, 593, 594, 595, 596, 597, 598,
       599, 600, 601, 602, 603, 604, 605, 606, 607, 608, 609, 610, 611,
       612, 613, 614, 615, 616, 617, 618, 619, 620, 621, 622, 623, 624,
       625, 626, 627, 628, 629, 630, 631, 632, 633, 634, 635, 636, 637,
       638, 639, 640, 641, 642, 643, 644, 645, 646, 647, 648, 649, 650,
       651, 652, 653, 654, 655, 656, 657, 658, 659, 660, 661, 662, 663,
       664, 665, 666, 667, 668, 669, 670, 671, 672, 673, 674, 675, 676,
       677, 678, 679, 680, 681, 682, 683, 684, 685, 686, 687, 688, 689,
       690, 691, 692, 693, 694, 695, 696, 697, 698, 699, 700, 701, 702,
       703, 704, 705, 706, 707, 708, 709, 710, 711, 712, 713, 714, 715,
       716, 717, 718, 719, 720, 721, 722, 723, 724, 725, 726, 727, 728,
       729, 730, 731, 732, 733, 734, 735, 736, 737, 738, 739, 740, 741,
       742, 743, 744, 745, 746, 747, 748, 749, 750, 751, 752, 753, 754,
       755, 756, 757, 758, 759, 760, 761, 762, 763, 764, 765, 766, 767,
       768, 769, 770, 771, 772, 773, 774, 775, 776, 777, 778, 779, 780,
       781, 782, 783, 784, 785, 786, 787, 788, 789, 790, 791, 792, 793,
       794, 795, 796, 797, 798, 799, 800, 801, 802, 803, 804, 805, 806,
       807, 808, 809, 810, 811, 812, 813, 814, 815, 816, 817, 818, 819,
       820, 821, 822, 823, 824, 825, 826, 827, 828, 829, 830, 831, 832,
       833, 834, 835, 836, 837, 838, 839, 840, 841, 842, 843, 844, 845,
       846, 847, 848, 849, 850, 851, 852, 853, 854, 855, 856, 857, 858,
       859, 860, 861, 862, 863, 864, 865, 866, 867, 868, 869, 870, 871,
       872, 873, 874, 875, 876, 877, 878, 879, 880, 881, 882, 883, 884,
       885, 886, 887, 888, 889, 890, 891], dtype=int64)
```

In [15]: `df["Survived"].unique()`

Out[15]: `array([0, 1], dtype=int64)`

In [16]: `df["Pclass"].unique()`

Out[16]: `array([3, 1, 2], dtype=int64)`

In [17]: `df["Parch"].unique()`

Out[17]: `array([0, 1, 2, 5, 3, 4, 6], dtype=int64)`

In [18]: `df["Embarked"].unique()`

Out[18]: `array(['S', 'C', 'Q', nan], dtype=object)`

In [19]: `df["SibSp"].unique()`

Out[19]: `array([1, 0, 3, 4, 2, 5, 8], dtype=int64)`

### 3)  Identify symmetric and asymmetric binary attributes from data sets and display all values.

In [ ]:

In [ ]:

In [ ]:
```
symmetric=["Sex"]
asymmetric=["Survived"]
```

### 4)  For each quantitative attribute, calculate its average, standard deviation, minimum, mode, range and maximum values.

In [43]:
```python
quantitative = ["PassengerId","Survived","Pclass","Age","SibSp","Parch","Fare"]
for i in quantitative :
    print(    )
    print(i)

    print("\tmean",df[i].mean())
    print("\tstandard deviation",df[i].std())
    print("\tminimum",df[i].min())
    print("\tmaximum",df[i].max())
    print("\tmode",df[i].mode()[0])
    print("\trange",df[i].max()-df[i].min())
```

```
PassengerId
        mean 446.0
        standard deviation 257.3538420152301
        minimum 1
        maximum 891
        mode 1
        range 890

Survived
        mean 0.3838383838383838
        standard deviation 0.4865924542648585
        minimum 0
        maximum 1
        mode 0
        range 1

Pclass
        mean 2.308641975308642
        standard deviation 0.8360712409770513
        minimum 1
        maximum 3
        mode 3
        range 2

Age
        mean 29.69911764705882
        standard deviation 14.526497332334044
        minimum 0.42
        maximum 80.0
        mode 24.0
        range 79.58

SibSp
        mean 0.5230078563411896
        standard deviation 1.1027434322934275
        minimum 0
        maximum 8
        mode 0
        range 8

Parch
        mean 0.38159371492704824
        standard deviation 0.8060572211299559
        minimum 0
        maximum 6
        mode 0
        range 6

Fare
        mean 32.204207968574636
        standard deviation 49.693428597180905
        minimum 0.0
        maximum 512.3292
        mode 8.05
        range 512.3292
```

**6) For the qualitative attribute (class), count the frequency for each of its distinct values.**

```
In [41]: df["Pclass"].value_counts()
```

```
Out[41]: 3    491
         1    216
         2    184
         Name: Pclass, dtype: int64
```

**7) It is also possible to display the summary for all the attributes simultaneously in a table using the describe() function. If an attribute is quantitative, it will display its mean, standard deviation and various quantiles (including minimum, median, and maximum) values. If an attribute is qualitative, it will display its number of unique values and the top (most frequent) values.**

```
In [48]: print(df.describe(include="all"))
```

```
              PassengerId    Survived      Pclass                      Name   Sex  \
count        891.000000  891.000000  891.000000                       891   891
unique              NaN         NaN         NaN                       891     2
top                 NaN         NaN         NaN  Braund, Mr. Owen Harris  male
freq                NaN         NaN         NaN                         1   577
mean         446.000000    0.383838    2.308642                       NaN   NaN
std          257.353842    0.486592    0.836071                       NaN   NaN
min            1.000000    0.000000    1.000000                       NaN   NaN
25%          223.500000    0.000000    2.000000                       NaN   NaN
50%          446.000000    0.000000    3.000000                       NaN   NaN
75%          668.500000    1.000000    3.000000                       NaN   NaN
max          891.000000    1.000000    3.000000                       NaN   NaN

                   Age       SibSp       Parch  Ticket        Fare    Cabin  \
count       714.000000  891.000000  891.000000     891  891.000000      204
unique             NaN         NaN         NaN     681         NaN      147
top                NaN         NaN         NaN  347082         NaN  B96 B98
freq               NaN         NaN         NaN       7         NaN        4
mean         29.699118    0.523008    0.381594     NaN   32.204208      NaN
std          14.526497    1.102743    0.806057     NaN   49.693429      NaN
min           0.420000    0.000000    0.000000     NaN    0.000000      NaN
25%          20.125000    0.000000    0.000000     NaN    7.910400      NaN
50%          28.000000    0.000000    0.000000     NaN   14.454200      NaN
75%          38.000000    1.000000    0.000000     NaN   31.000000      NaN
max          80.000000    8.000000    6.000000     NaN  512.329200      NaN

        Embarked
count        889
unique         3
top            S
freq         644
mean         NaN
std          NaN
min          NaN
25%          NaN
50%          NaN
75%          NaN
max          NaN
```

```
In [49]: df.describe()
```

Out[49]:

|       | PassengerId | Survived   | Pclass     | Age        | SibSp      | Parch      | Fare       |
|-------|-------------|------------|------------|------------|------------|------------|------------|
| count | 891.000000  | 891.000000 | 891.000000 | 714.000000 | 891.000000 | 891.000000 | 891.000000 |
| mean  | 446.000000  | 0.383838   | 2.308642   | 29.699118  | 0.523008   | 0.381594   | 32.204208  |
| std   | 257.353842  | 0.486592   | 0.836071   | 14.526497  | 1.102743   | 0.806057   | 49.693429  |
| min   | 1.000000    | 0.000000   | 1.000000   | 0.420000   | 0.000000   | 0.000000   | 0.000000   |
| 25%   | 223.500000  | 0.000000   | 2.000000   | 20.125000  | 0.000000   | 0.000000   | 7.910400   |
| 50%   | 446.000000  | 0.000000   | 3.000000   | 28.000000  | 0.000000   | 0.000000   | 14.454200  |
| 75%   | 668.500000  | 1.000000   | 3.000000   | 38.000000  | 1.000000   | 0.000000   | 31.000000  |
| max   | 891.000000  | 1.000000   | 3.000000   | 80.000000  | 8.000000   | 6.000000   | 512.329200 |

In [50]: `df.describe(include = "object")`

Out[50]:

|  | Name | Sex | Ticket | Cabin | Embarked |
|---|---|---|---|---|---|
| **count** |  | 891 | 891 | 891 | 204 | 889 |
| **unique** |  | 891 | 2 | 681 | 147 | 3 |
| **top** | Braund, Mr. Owen Harris | male | 347082 | B96 B98 | S |
| **freq** |  | 1 | 577 | 7 | 4 | 644 |

## 8) For multivariate statistics, you can compute the covariance and correlation between pairs of attributes.

In [54]: `df.corr()`

```
C:\Users\student\AppData\Local\Temp\ipykernel_2196\1134722465.py:1: FutureWarning: The default value of n
umeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only v
alid columns or specify the value of numeric_only to silence this warning.
  df.corr()
```

Out[54]:

|  | PassengerId | Survived | Pclass | Age | SibSp | Parch | Fare |
|---|---|---|---|---|---|---|---|
| **PassengerId** | 1.000000 | -0.005007 | -0.035144 | 0.036847 | -0.057527 | -0.001652 | 0.012658 |
| **Survived** | -0.005007 | 1.000000 | -0.338481 | -0.077221 | -0.035322 | 0.081629 | 0.257307 |
| **Pclass** | -0.035144 | -0.338481 | 1.000000 | -0.369226 | 0.083081 | 0.018443 | -0.549500 |
| **Age** | 0.036847 | -0.077221 | -0.369226 | 1.000000 | -0.308247 | -0.189119 | 0.096067 |
| **SibSp** | -0.057527 | -0.035322 | 0.083081 | -0.308247 | 1.000000 | 0.414838 | 0.159651 |
| **Parch** | -0.001652 | 0.081629 | 0.018443 | -0.189119 | 0.414838 | 1.000000 | 0.216225 |
| **Fare** | 0.012658 | 0.257307 | -0.549500 | 0.096067 | 0.159651 | 0.216225 | 1.000000 |

In [55]: `df.cov()`

```
C:\Users\student\AppData\Local\Temp\ipykernel_2196\1545644723.py:1: FutureWarning: The default value of n
umeric_only in DataFrame.cov is deprecated. In a future version, it will default to False. Select only va
lid columns or specify the value of numeric_only to silence this warning.
  df.cov()
```

Out[55]:

|  | PassengerId | Survived | Pclass | Age | SibSp | Parch | Fare |
|---|---|---|---|---|---|---|---|
| **PassengerId** | 66231.000000 | -0.626966 | -7.561798 | 138.696504 | -16.325843 | -0.342697 | 161.883369 |
| **Survived** | -0.626966 | 0.236772 | -0.137703 | -0.551296 | -0.018954 | 0.032017 | 6.221787 |
| **Pclass** | -7.561798 | -0.137703 | 0.699015 | -4.496004 | 0.076599 | 0.012429 | -22.830196 |
| **Age** | 138.696504 | -0.551296 | -4.496004 | 211.019125 | -4.163334 | -2.344191 | 73.849030 |
| **SibSp** | -16.325843 | -0.018954 | 0.076599 | -4.163334 | 1.216043 | 0.368739 | 8.748734 |
| **Parch** | -0.342697 | 0.032017 | 0.012429 | -2.344191 | 0.368739 | 0.649728 | 8.661052 |
| **Fare** | 161.883369 | 6.221787 | -22.830196 | 73.849030 | 8.748734 | 8.661052 | 2469.436846 |

In [11]:

Out[11]:

|  | PassengerId | Survived | Pclass | Age | SibSp | Parch | Fare |
|---|---|---|---|---|---|---|---|
| **PassengerId** | 1.000000 | -0.005007 | -0.035144 | 0.036847 | -0.057527 | -0.001652 | 0.012658 |
| **Survived** | -0.005007 | 1.000000 | -0.338481 | -0.077221 | -0.035322 | 0.081629 | 0.257307 |
| **Pclass** | -0.035144 | -0.338481 | 1.000000 | -0.369226 | 0.083081 | 0.018443 | -0.549500 |
| **Age** | 0.036847 | -0.077221 | -0.369226 | 1.000000 | -0.308247 | -0.189119 | 0.096067 |
| **SibSp** | -0.057527 | -0.035322 | 0.083081 | -0.308247 | 1.000000 | 0.414838 | 0.159651 |
| **Parch** | -0.001652 | 0.081629 | 0.018443 | -0.189119 | 0.414838 | 1.000000 | 0.216225 |
| **Fare** | 0.012658 | 0.257307 | -0.549500 | 0.096067 | 0.159651 | 0.216225 | 1.000000 |

**9) Display the histogram for Age attribute by discretizing it into 8 separate bins and counting the frequency for each bin.**
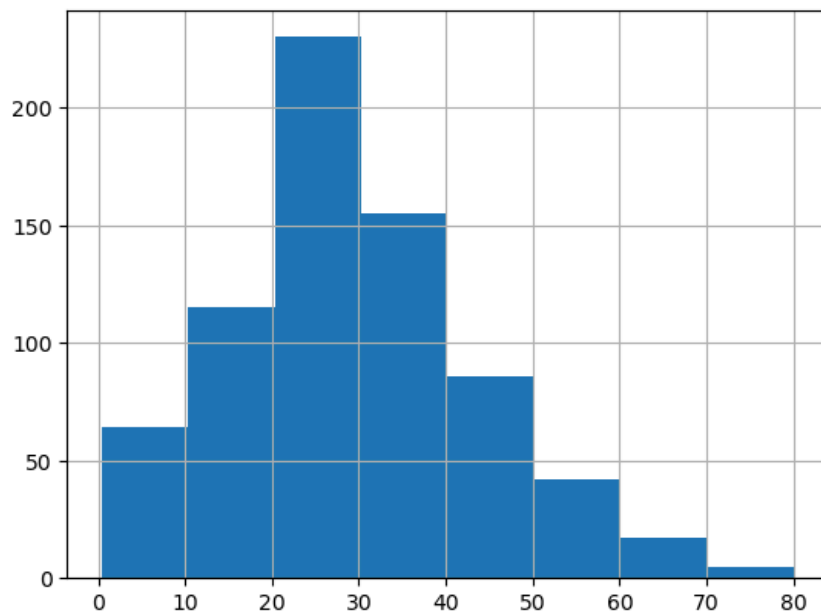
In [60]: 
```python
import matplotlib.pyplot as plt
```

In [64]: 
```python
plt.hist(df["Age"],bins =8)
plt.grid()
```
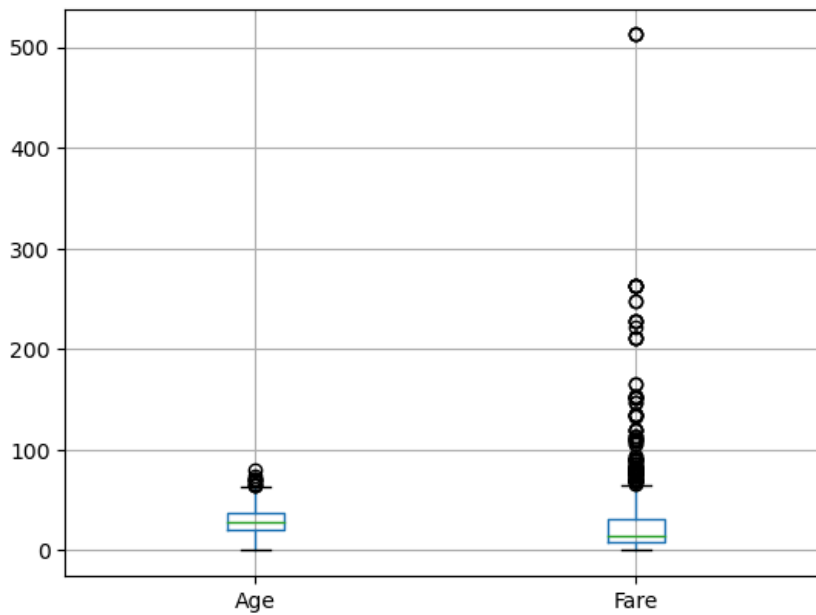


In [74]: 
```python
df["Age"].hist(bins=8)
```

Out[74]: <Axes: >

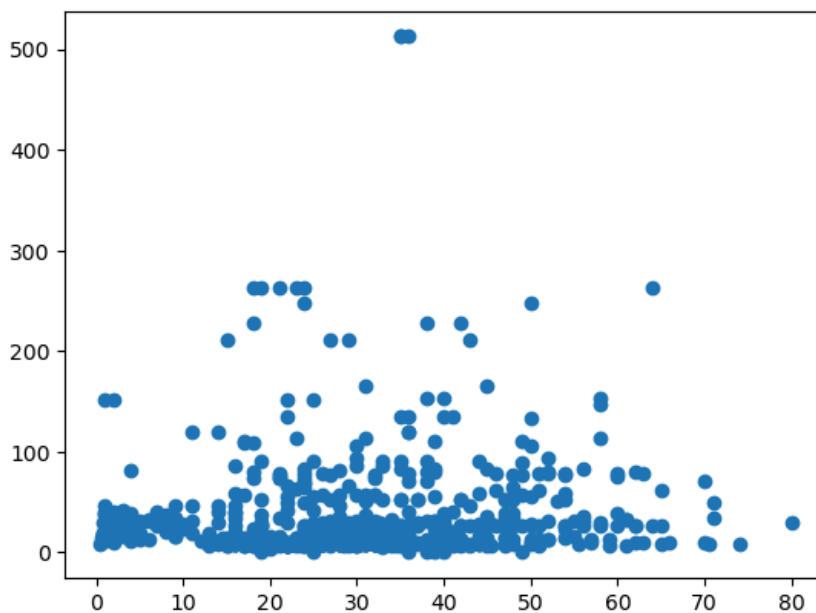**10) A boxplot can also be used to show the distribution of values for each attribute.**

In [73]: `df.boxplot(["Age","Fare"])`

Out[73]: `<Axes: >`



**11) Display scatter plot for any 5 pair of attributes , we can use a scatter plot to visualize their joint distribution.**

In [75]:
```
#onepare
plt.scatter(df["Age"],df["Fare"])
```

Out[75]: `<matplotlib.collections.PathCollection at 0x1ae13ad9840>`



In [ ]:

In [ ]: