





The dataset is created by IBM data scientists to uncover the factors that lead to employee attrition and explore important questions such as 'show me a breakdown of distance from home by job role and attrition' or 'compare average monthly income by education and attrition'.

The name of the dataset attributes are self explanatory. Some required attributes are explained below.

- Education
  - Below College
  - College
  - Bachelor
  - Master
  - Doctor
- EnvironmentSatisfaction
  - Low
  - Medium
  - High
  - Very High
- JobInvolvement
  - Low
  - Medium
  - High
  - Very High
- JobSatisfaction
  - Low
  - Medium
  - High
  - Very High
- PerformanceRating
  - Low
  - Good
  - Excellent
  - Outstanding
- RelationshipSatisfaction
  - Low
  - Medium
  - High
  - Very High
- WorkLifeBalance
  - Bad
  - Good
  - Better
  - Best

There are few libraries required to perform encoding variables:

- pandas - It helps to retrieve datasets, handle missing data and in data wrangling.
- Numpy - It helps to perform numerical operations in the dataset.
- seaborn and matplotlib- It helps to data visualization.
- warnings - It helps to segregate the warnings.

### Import libraries

```
#Select the cell and click on run icon
import numpy as np
import pandas as pd
import seaborn as sns

import matplotlib.pyplot as plt
import matplotlib
import warnings
warnings.filterwarnings('ignore')
```

#### Instruction:

Download the **HR-Employee-Attrition.csv** dataset file from Course Resources and upload it in the lab using Up arrow shown below

View Tab

### Read the dataset

#### Note 16:

- The **df** is dataframe to store the data imported from the csv as rows and columns table format.
- The **head()** function helps to view first few data present in the **df** dataframe.

```
#Select the cell and click on run icon
df = pd.read_csv("HR-Employee-Attrition.csv")
df.head()
```

```
dtype: int64
```

**Observations:**

There are no null values present in the `df` dataframe.

## Now, print the count of values for each column

```
for i in df.columns:
    print(i, ":", df[i].value_counts())
    print("-" * 40)
    print("-" * 40)
```

Age	:	35	78
		34	77
		36	69
		31	69
		29	68
		32	61
		30	60
		33	58
		38	58
		40	57
		37	50
		77	48

5 rows x 35 columns

#### Note 17:

The **columns** provides the name of the columns present in the **df** dataframe.

```
#Select the cell and click on run icon
df.columns
```

```
Index(['Age', 'Attrition', 'BusinessTravel', 'DailyRate', 'Department',
       'DistanceFromHome', 'Education', 'EducationField', 'EmployeeCount',
       'EmployeeNumber', 'EnvironmentSatisfaction', 'Gender', 'HourlyRate',
       'JobInvolvement', 'JobLevel', 'JobRole', 'JobSatisfaction', 'NumCompaniesWorked',
       'MaritalStatus', 'MonthlyIncome', 'MonthlyRate', 'NumCompaniesWorked',
       'Over18', 'OverTime', 'PercentSalaryHike', 'PerformanceRating',
       'RelationshipSatisfaction', 'StandardHours', 'StockOptionLevel',
       'TotalWorkingYears', 'TrainingTimesLastYear', 'WorkLifeBalance',
       'YearsAtCompany', 'YearsInCurrentRole', 'YearsSinceLastPromotion',
       'YearsWithCurrManager'],
      dtype='object')
```

#### Observations we can draw from the above output:

- It denotes the name of the columns present in the **df** dataframe.

#### Note 18:

The **shape** represents the number of elements in each dimension and returns a tuple with each index having the number of corresponding elements.

```
#Select the cell and click on run icon
df.shape
```

(1470, 35)

#### Observations we can draw from the above output:

- There are 1470 rows and 35 columns present in the **df** dataframe.

### Check if there are any null values

```
#Select the cell and click on run icon
df.isna().sum()
```

```
Age      0
Attrition 0
BusinessTravel 0
DailyRate 0
Department 0
DistanceFromHome 0
Education 0
EducationField 0
EmployeeCount 0
EmployeeNumber 0
EnvironmentSatisfaction 0
Gender 0
HourlyRate 0
JobInvolvement 0
JobLevel 0
JobRole 0
JobSatisfaction 0
MaritalStatus 0
MonthlyIncome 0
MonthlyRate 0
NumCompaniesWorked 0
Over18 0
OverTime 0
PercentSalaryHike 0
PerformanceRating 0
RelationshipSatisfaction 0
StandardHours 0
StockOptionLevel 0
TotalWorkingYears 0
TrainingTimesLastYear 0
WorkLifeBalance 0
YearsAtCompany 0
YearsInCurrentRole 0
YearsSinceLastPromotion 0
YearsWithCurrManager 0
dtype: int64
```

#### Observations:

- There are no null values present in the **df** dataframe.

### Now, print the count of values for each column

```
for i in df.columns:
    print(i, ":", df[i].value_counts())
    print(" " * 40)
    print(" " * 40)
```

```
Age : 35    78
      34    77
      36    69
      31    69
      32    68
      32    61
      30    60
      33    58
      38    58
      46    57
      37    50
      27    48
      28    48
      42    46
      39    42
      45    41
      41    40
      26    39
      44    33
      46    33
      43    32
      50    30
      25    26
      24    26
      49    24
      47    22
      55    22
      51    19
      53    19
      48    19
      54    18
      52    18
      22    16
      56    14
      23    14
      58    14
      21    13
      20    11
      59    10
      19    9
      18    8
      60    5
      57    4
Name: Age, dtype: int64

Attrition : No      1233
           Yes       237
Name: Attrition, dtype: int64

BusinessTravel : Travel_Rarely      1043
                Travel_Frequently    277
                Non-Travel           130
Name: BusinessTravel, dtype: int64

DailyRate : 691      6
            408      5
            530      5
            1329      5
            1082      5
            650      1
            279      1
            316      1
            314      1
            628      1
Name: DailyRate, Length: 886, dtype: int64

Department : Research & Development      961
            Sales                        446
            Human Resources                63
Name: Department, dtype: int64

DistanceFromHome : 2      211
                  1      208
                  10     86
                  9     85
                  3     84
                  7     84
                  8     80
                  5     65
                  4     64
                  6     59
                  16     32
                  11     29
                  24     28
                  23     27
                  29     27
                  15     26
                  18     26
                  26     25
                  25     25
                  20     25
                  28     23
                  19     22
                  14     21
                  12     20
                  17     20
                  22     19
                  13     19
                  21     18
                  27     12
Name: DistanceFromHome, dtype: int64

Education : 3      572
            4      398
            2      282
            1      170
            5      48
Name: Education, dtype: int64

EducationField : Life Sciences      606
                Medical             464
                Marketing            159
                Technical Degree     132
                Other                82
                Human Resources      27
Name: EducationField, dtype: int64

EmployeeCount : 1      1470
Name: EmployeeCount, dtype: int64

EmployeeNumber : 1      1
                1391      1
                1389      1
                1387      1
                1383      1
                659      1
                657      1
                656      1
                655      1
                2068      1
Name: EmployeeNumber, Length: 1470, dtype: int64

EnvironmentSatisfaction : 3      453
                          4      446
                          2      297
                          1      284
Name: EnvironmentSatisfaction, dtype: int64

Gender : Male      882
        Female     588
Name: Gender, dtype: int64

HourlyRate : 66      29
            88      28
            42      28
            48      28
            64      28
            ..
            31      15
            53      14
            68      14
            38      13
            34      12
Name: HourlyRate, Length: 71, dtype: int64

JobInvolvement : 3      868
                 2      375
                 4      144
                 1      83
Name: JobInvolvement, dtype: int64

JobLevel : 1      543
           2      534
           3      218
           4      106
           5      69
Name: JobLevel, dtype: int64

JobRole : Sales Executive      326
         Research Scientist    292
         Laboratory Technician 259
         Manufacturing Director 145
         Healthcare Representative 131
         Manager               102
         Sales Representative    83
         Research Director       80
         Human Resources         52
Name: JobRole, dtype: int64

JobSatisfaction : 4      459
                 3      442
                 1      289
                 2      280
Name: JobSatisfaction, dtype: int64

MaritalStatus : Married      673
               Single        470
               Divorced       527
Name: MaritalStatus, dtype: int64

MonthlyIncome : 2342      4
              6544      1
              2741      3
              2559      3
              2610      3
              7104      1
              2773      1
              19513      1
              3447      1
              4404      1
Name: MonthlyIncome, Length: 1349, dtype: int64

MonthlyRate : 4223      3
             9150      3
             9558      2
             12618      2
             22074      2
             14561      1
             2671      1
             5715      1
             11757      1
             10228      1
Name: MonthlyRate, Length: 1427, dtype: int64

NumCompaniesWorked : 1      521
                    0      197
                    3      159
                    2      146
                    4      139
                    7      74
                    6      70
                    5      63
                    9      52
                    8      49
Name: NumCompaniesWorked, dtype: int64

Over18 : Y      1470
Name: Over18, dtype: int64

OverTime : No      1054
          Yes       416
Name: OverTime, dtype: int64

PercentSalaryHike : 11      210
                   13      199
                   14      201
                   12      198
                   15      191
                   18      89
                   17      82
                   16      78
                   19      76
                   22      56
                   20      55
                   21      48
                   23      28
                   24      21
                   25      18
Name: PercentSalaryHike, dtype: int64

PerformanceRating : 3      1244
                   4      126
Name: PerformanceRating, dtype: int64

RelationshipSatisfaction : 3      459
                          4      432
                          2      303
                          1      276
Name: RelationshipSatisfaction, dtype: int64

StandardHours : 80      1470
Name: StandardHours, dtype: int64

StockOptionLevel : 0      631
                  1      596
                  2      158
                  3      85
Name: StockOptionLevel, dtype: int64

TotalWorkingYears : 10      202
                   6      197
                   8      103
                   9      96
                   5      88
                   7      81
                   1      81
                   4      63
                   12      48
                   3      42
                   15      40
                   16      37
                   11      36
                   13      36
                   21      34
                   17      33
                   2      31
                   14      31
                   20      30
                   18      27
                   19      22
                   23      22
                   22      21
                   24      18
                   25      14
                   28      14
                   26      14
                   9      11
                   31      9
                   32      9
                   30      7
                   33      7
                   27      6
                   36      6
                   34      5
                   37      4
                   35      3
                   40      2
                   38      2
Name: TotalWorkingYears, dtype: int64

TrainingTimesLastYear : 2      547
                       3      491
                       4      123
                       5      119
                       1      71
                       6      65
                       1      64
Name: TrainingTimesLastYear, dtype: int64

WorkLifeBalance : 3      893
                 4      844
                 4      153
                 2      82
Name: WorkLifeBalance, dtype: int64

YearsAtCompany : 5      196
                3      128
                2      127
                10     120
                4      110
                9      90
                8      80
                6      78
                11     32
                20     27
                13     24
                15     20
                14     18
                22     15
                21     14
                16     12
                19     11
                17     9
                24     6
                33     5
                32     4
                13     4
                26     4
                31     3
                27     2
                36     2
                23     2
                37     1
                40     1
                34     1
                30     1
Name: YearsAtCompany, dtype: int64

YearsInCurrentRole : 2      372
                    3      135
                    4      104
                    7      222
                    5      187
                    1      171
                    6      37
                    5      36
                    11     22
                    13     14
                    14     14
                    17     7
                    15     5
                    14     5
                    16     2
Name: YearsInCurrentRole, dtype: int64

YearsSinceLastPromotion : 0      581
                        2      159
                        4      76
                        3      52
                        5      45
                        6      32
                        11     24
                        8      18
                        9      17
                        15     13
                        13     10
                        12     10
                        14     9
                        10     8
Name: YearsSinceLastPromotion, dtype: int64

YearsWithCurrManager : 2      344
                      3      163
                      7      216
                      3      142
                      8      107
                      4      98
                      1      76
                      6      29
                      5      31
                      11     22
                      12     18
                      13     14
                      17     7
                      15     5
                      14     5
                      16     2
Name: YearsWithCurrManager, dtype: int64
```

#### Observations we can draw from the above output:

- It is an imbalanced data as value count in **attrition** column is imbalanced.

Over18 is Y (Yes) across all the employees and it is not beneficial for further use.

StandardHours is 80 which is common for all and it is not beneficial for further use.

EmployeeCount, EmployeeNumber are unique which are also not going to help us predict our end result.

### Now, print the unique values in each column to understand the above output more clearly

```
for i in df.columns:
    print(i, ":", df[i].unique())
    print(" " * 40)
    print(" " * 40)
```



Age : [41 49 37 33 27 32 59 30 38 36 35 29 31 34 28 22 53 24 21 42 44 46 39 43  
50 26 48 55 43 56 51 52 40 54 56 52 25 19 57 52 47 18 60]

Attrition : ['Yes' 'No']

BusinessTravel : ['Travel\_Rarely' 'Travel\_Frequently' 'Non-Travel']

DailyRate : [1102 279 1373 1392 591 1005 1324 1358 216 1299 809 153 670 1346  
103 1389 334 1123 1219 371 673 1218 419 391 699 1282 1125 691  
477 705 924 1459 125 895 813 1273 869 890 852 1141 464 1240  
1357 894 721 1360 1065 408 2121 1229 626 1434 1448 1097 1443 515  
853 1142 655 1115 427 653 989 1435 1223 836 1195 1339 664 318  
1225 1328 1082 548 132 746 776 193 397 945 1214 111 573 1153  
1400 541 432 289 669 530 632 1234 639 1093 1271 1253 120 682  
489 807 827 871 665 1040 1420 240 1280 534 1456 658 142 1127  
1031 1189 1354 1467 922 394 1312 750 441 684 249 841 147 528  
594 470 997 542 802 1355 1150 1229 959 1033 1316 364 438 689  
201 1427 857 933 1181 1395 662 1436 194 967 1496 1169 1145 630  
303 1256 440 1450 1452 465 702 1157 602 1480 1268 713 134 526  
1380 140 629 1356 128 1084 931 692 1063 913 894 556 1344 290  
138 926 1261 472 1002 878 905 1180 121 1136 635 1151 644 1045  
829 1242 1469 896 992 1052 1147 1396 663 119 979 319 1413 944  
1323 532 818 854 1034 771 1401 1431 976 1411 1300 252 1327 832  
1017 1159 504 505 916 1247 685 269 1416 833 307 1311 128 488  
529 1210 1463 675 1385 1403 452 666 1158 228 996 728 1315 322  
1479 797 1070 442 496 1372 920 688 1449 1117 636 506 444 950  
889 555 230 1232 566 1302 812 1476 218 1132 1105 906 849 390  
106 1249 192 553 117 185 1091 723 1180 1202 888 1377 1018 1275 798  
672 1162 508 1482 559 210 928 1001 549 1124 758 570 1130 1192  
343 144 1296 1309 483 810 544 1062 1319 641 1332 756 845 593  
1171 350 921 1144 143 1046 575 136 1283 755 304 1178 329 1362  
1371 202 253 164 1107 759 1025 982 821 186 480 1473 891 1063  
645 1490 317 422 1485 1368 1446 296 1398 1349 986 1099 116 1499  
983 1009 1303 1274 1277 587 413 1276 988 1474 163 267 619 302  
443 828 561 422 1232 1306 1094 509 735 806 817 1410 207 1442 493  
535 1495 446 1245 703 823 1246 622 1287 408 254 1365 538 525  
558 782 362 1236 1112 204 1343 604 1216 646 160 238 1397 306  
991 482 1176 912 876 727 376 571 384 791 1211 1243 1032 1225  
529 662 508 680 970 1179 294 314 316 654 148 381 217 501  
650 141 804 973 1090 346 430 268 167 621 527 883 954 310  
725 715 657 1146 162 376 101 384 791 1211 1243 1032 1225  
805 213 118 676 1252 286 1258 932 1041 859 720 946 1184 436  
589 760 887 1318 625 180 586 1012 661 930 342 1230 1271 1278  
1138 130 300 582 1818 1269 375 385 1263 1222 341 868 1221 102  
881 1383 1075 374 1086 781 177 550 1425 1454 617 1085 995 1122  
618 546 462 1198 1272 154 1137 1188 189 1333 867 263 938 129  
616 498 1404 1053 289 1376 231 152 882 903 1379 335 722 461  
974 1126 840 1134 248 955 939 1391 1206 287 141 109 1066 277  
466 1055 265 135 247 1035 266 145 1038 1234 1109 1089 788 124  
650 1186 1464 796 415 769 1003 1366 330 1492 1204 309 1330 469  
697 1262 1050 770 406 203 1308 984 439 793 1451 1182 1174 490  
718 433 773 603 874 367 199 481 647 1384 902 819 862 1457  
977 942 1402 1421 1361 917 200 150 179 696 116 363 107 1465  
458 1212 1103 966 1010 326 1098 969 1167 694 1320 536 373 599  
251 131 237 1429 648 735 531 429 968 879 640 412 848 360  
1138 325 1322 239 1030 634 924 236 1060 935 495 282 206 943  
523 507 601 855 1291 1405 1369 999 1202 285 404 736 1498 1200  
1439 499 205 683 1462 949 652 332 1475 337 971 1174 667 560  
1372 1285 1205 359 403 377 592 1445 1221 866 981 447 1316 749  
990 405 115 790 830 1193 1423 467 271 410 1083 516 224 136  
1029 333 1440 674 1342 898 824 492 598 740 888 1288 104 1108  
479 1351 474 617 884 1370 264 1056 561 176 897 970 1054 428 191  
1387 170 208 671 731 737 1470 365 763 867 486 772 301 311  
584 880 392 148 708 1259 786 370 678 146 581 918 1238 585  
741 552 368 711 543 964 794 621 176 897 970 1054 428 191  
211 1079 590 305 953 478 1375 244 511 1294 196 734 1239 1253  
1128 1336 234 766 261 1194 431 572 1422 1297 574 355 207 706  
780 736 414 352 1324 459 2154 1131 835 1372 1266 783 219 1213  
1096 1251 1394 605 1064 1337 937 157 754 1168 155 1444 189 911  
1321 1154 557 642 801 161 1382 1037 105 582 704 345 1120 1378  
468 613 1023 628]

Department : ['Sales' 'Research & Development' 'Human Resources']

DistanceFromHome : [ 1 8 2 3 24 23 27 16 15 26 19 21 5 11 9 7 6 10 4 25 12 18 29 22  
14 20 28 17 13]

Education : [2 1 4 3 5]

EducationField : ['Life Sciences' 'Other' 'Medical' 'Marketing' 'Technical Degree'  
'Human Resources']

EmployeeCount : [1]

EmployeeNumber : [ 1 2 4 ... 2064 2065 2068]

EnvironmentSatisfaction : [2 3 4 1]

Gender : ['Female' 'Male']

HourlyRate : [ 94 61 92 56 40 79 81 67 44 84 49 31 93 50 51 80 96 78  
45 82 83 83 58 12 48 42 4 86 97 78 37 73 98 36 47  
71 30 43 99 59 95 57 76 87 66 55 32 52 70 62 64 63 60  
100 46 39 77 35 91 54 34 90 65 88 85 89 68 69 74 38]

JobInvolvement : [3 2 4 1]

JobLevel : [2 1 3 4 5]

JobRole : ['Sales\_Executive' 'Research\_Scientist' 'Laboratory\_Technician'  
'Manufacturing\_Director' 'Healthcare\_Representative' 'Manager'  
'Sales\_Representative' 'Research\_Director' 'Human\_Resources']

JobSatisfaction : [4 2 3 1]

MaritalStatus : ['Single' 'Married' 'Divorced']

MonthlyIncome : [5993 5130 2090 ... 9991 5390 4404]

MonthlyRate : [19479 24907 2396 ... 5174 13243 10228]

NumCompaniesWorked : [8 1 6 9 0 4 5 2 7 3]

Over18 : ['Y']

OverTime : ['Yes' 'No']

PercentSalaryHike : [11 23 15 12 13 20 22 21 17 14 16 18 19 24 25]

PerformanceRating : [3 4]

RelationshipSatisfaction : [1 4 2 3]

StandardHours : [80]

StockOptionLevel : [0 1 3 2]

TotalWorkingYears : [8 10 7 6 12 1 17 5 3 21 13 0 26 24 22 9 19 2 23 14 15 4 29 28  
21 25 20 11 16 37 30 40 18 36 34 32 33 35 2]

TrainingTimesLastYear : [0 3 2 5 1 4 6]

WorkLifeBalance : [1 3 2 4]

YearsAtCompany : [ 6 10 0 8 2 7 1 9 5 4 25 3 12 14 22 15 27 21 17 11 13 37 16 20  
40 24 33 19 36 18 29 31 32 34 26 30 23]

YearsInCurrentRole : [ 4 7 0 2 5 9 8 3 6 13 3 15 14 16 11 10 12 18 17]

YearsSinceLastPromotion : [ 0 1 3 2 7 4 8 6 5 15 9 13 12 10 11 14]

YearsWithCurrManager : [ 5 7 0 2 6 8 3 11 17 1 4 12 9 10 15 13 16 14]

Note 19:

The **drop()** function helps to drop the columns which are not important features in the dataset.

In [28]: df=df.drop(['Over18','EmployeeNumber','EmployeeCount','StandardHours'],axis=1)

Note 20:

After dropping, check the number of columns and rows present in the dataset.

In [29]: df.shape

Out[29]: (1470, 31)

Note 21:

The **columns** provides the name of the columns present in the **df** dataframe.

In [30]: df.columns

Out[30]: Index(['Age', 'Attrition', 'BusinessTravel', 'DailyRate', 'Department',  
'DistanceFromHome', 'Education', 'EducationField', 'EnvironmentSatisfaction', 'Gender', 'HourlyRate', 'JobInvolvement',  
'JobLevel', 'JobRole', 'JobSatisfaction', 'MaritalStatus',  
'MonthlyIncome', 'MonthlyRate', 'NumCompaniesWorked', 'OverTime',  
'PercentSalaryHike', 'PerformanceRating', 'RelationshipSatisfaction',  
'StockOptionLevel', 'TotalWorkingYears', 'TrainingTimesLastYear',  
'WorkLifeBalance', 'YearsAtCompany', 'YearsInCurrentRole', 'YearsSinceLastPromotion',  
'YearsWithCurrManager', 'YearsInCurrentRole',  
'YearsSinceLastPromotion', 'YearsWithCurrManager'],  
dtype='object')

Observations from the above output:

We can see a lot of columns which deal with satisfaction, let's observe them closely.

In [31]: df[['RelationshipSatisfaction','JobSatisfaction','EnvironmentSatisfaction','JobInvolvement']]

Out[31]:

	RelationshipSatisfaction	JobSatisfaction	EnvironmentSatisfaction	JobInvolvement
0	1	4	2	3
1	4	2	3	2
2	2	3	4	2
3	3	3	4	3
4	507	601	2	1
...	...	...	...	...
1465	3	4	3	4
1466	1	1	4	2
1467	2	2	2	4
1468	4	2	4	2
1469	1	3	2	4

1470 rows x 4 columns

Observations from the above output:

All the satisfactions are measured from 1 to 4 in incremental order, along with WorkLifeBalance.

Higher the value, higher the satisfaction.

So, we can combine all the columns that convey satisfaction detail and make them one.

Note 22:

Let us now calculate the mean of all the types of satisfaction. Based on a condition if the mean value is greater than 2.35 then it returns one else zero

In [32]: df['TotalSatisfaction\_mean'] = (df['RelationshipSatisfaction'] + df['EnvironmentSatisfaction']  
+ df['JobSatisfaction'] + df['JobInvolvement'] + df['WorkLifeBalance'])/5

def Satif(df):  
if df['TotalSatisfaction\_mean'] > 2.35 :  
return 1  
else:  
return 0

df['Satif'] = df.apply(lambda df:Satif(df) ,axis = 1)  
df['Satif']

Out[32]:

0	0
1	1
2	1
3	1
4	1
...	...
1465	1
1466	0
1467	1
1468	1
1469	1

Name: Satif, Length: 1470, dtype: int64

In [33]: df.shape

Out[33]: (1470, 33)

Let's keep all the columns for now.

Note 23:

Let us now create a separate column for job satisfaction.

In [34]: df['JobSatif\_mean'] = (df['JobSatisfaction'] + df['JobInvolvement']) / 2

In [35]: df.shape

Out[35]: (1470, 34)

Note 24:

We can understand people who switch companies frequently have more tendency to leave a company. So let's create a new column called 'MovingPeople'.

In [36]: def MovingPeople(df):  
if df['NumCompaniesWorked'] > 4:  
return 1  
else:  
return 0  
df['MovingPeople'] = df.apply(lambda df:MovingPeople(df), axis = 1)  
df['MovingPeople']

Out[36]:

0	1
1	1
2	1
3	0
4	1
...	...
1465	0
1466	0
1467	0
1468	0
1469	0

Name: MovingPeople, Length: 1470, dtype: int64

In [37]: df.shape

Out[37]: (1470, 35)

Note 25:

Create a column using 'DistanceFromHome' column

In [38]: def LongDis(df):  
if df['DistanceFromHome'] > 11:  
return 1  
else:  
return 0  
df['LongDis'] = df.apply(lambda df:LongDis(df), axis = 1)  
df['LongDis']

Out[38]:

0	0
1	0
2	0
3	0
4	0
...	...
1465	1
1466	0
1467	0
1468	0
1469	0

Name: LongDis, Length: 1470, dtype: int64

Note 26:

Create a column using 'TrainingTimesLastYear' column

In [39]: def MiddleTraining(df):  
if df['TrainingTimesLastYear'] >= 3 and df['TrainingTimesLastYear'] <= 6:  
return 1  
else:  
return 0  
df['MiddleTraining'] = df.apply(lambda df:MiddleTraining(df), axis = 1)  
df['MiddleTraining']

Note 27:

Create a column to view number of years worked in each company

In [40]: df['Time\_in\_each\_comp'] = (df['Age'] - 20) / ((df['NumCompaniesWorked'] + 1)  
df['Time\_in\_each\_comp'])

Out[40]:

0	2.333333
1	1.500000
2	2.428571
3	6.500000
4	0.700000
...	...
1465	3.200000
1466	3.800000
1467	3.500000
1468	9.666667
1469	6.666667

Name: Time\_in\_each\_comp, Length: 1470, dtype: float64

In [41]: df.shape

Out[41]: (1470, 38)

Note 28:

Let us now understand the columns that are categorical and numerical in the **df** dataframe.

In [42]: numeric\_df=df.select\_dtypes(include=[np.number])  
categoric\_df=df.select\_dtypes(exclude=[np.number])

In [43]: numericcol=numeric\_df.columns.tolist()  
categorycol=categoric\_df.columns.tolist()

print ("Category :",categorycol)  
print ("N Numeric :",numericcol)

Category : ['Attrition', 'BusinessTravel', 'Department', 'EducationField', 'Gender', 'JobRole', 'MaritalStatus', 'OverTime']  
Numeric : ['Age', 'DailyRate', 'DistanceFromHome', 'Education', 'EnvironmentSatisfaction', 'HourlyRate', 'JobInvolvement', 'JobLevel', 'JobRole', 'JobSatisfaction', 'MaritalStatus', 'MonthlyIncome', 'MonthlyRate', 'NumCompaniesWorked', 'PercentSalaryHike', 'PerformanceRating', 'RelationshipSatisfaction', 'StockOptionLevel', 'TotalWorkingYears', 'TrainingTimesLastYear', 'WorkLifeBalance', 'YearsAtCompany', 'YearsInCurrentRole', 'YearsSinceLastPromotion', 'YearsWithCurrManager', 'TotalSatisfaction\_mean', 'Satif', 'JobSatif\_mean', 'MovingPeople', 'LongDis', 'MiddleTraining', 'BusinessTravel\_Travel\_Rarely', 'Department\_Research & Development', 'Department\_Sales', 'EducationField\_Life Sciences', 'EducationField\_Marketing', 'EducationField\_Medical', 'EducationField\_Other', 'EducationField\_Technical Degree', 'Gender\_Male', 'JobRole\_Human Resources', 'JobRole\_Laboratory Technician', 'JobRole\_Management', 'JobRole\_Manufacturing Director', 'JobRole\_Research Director', 'JobRole\_Research Scientist', 'JobRole\_Sales Executive', 'JobRole\_Sales Representative', 'MaritalStatus\_Married', 'MaritalStatus\_Single', 'OverTime\_Yes'], dtype='object'

In [47]: data.head()

Out[47]:

	Age	Education	JobLevel	MonthlyIncome	MonthlyRate	PercentSalaryHike	PerformanceRating	StockOptionLevel	TotalWorkingYears	Work
0	41	2	2	5993	19479	11	3	0	8	
1	49	1	2	5130	24907	23	4	1	10	
2	37	2	1	2090	2396	15	3	0	7	
3	33	4	1	2909	23159	11	3	0	8	
4	27	1	1	3468	16632	12	3	1	6	

5 rows x 43 columns

Observations from the above output:

We can see that all the variables present in the dataframe are in the numerical format.

Note: The rest of the encoding methods will be used in further lessons

The further operations come under feature Selection

Note: In this lesson, we saw the use of the feature engineering methods, but in the next lesson we are going to use one of these methods as a sub component for 'Exploratory Data Analysis'.

Powered by simplilearn

In [ ]: