

```
In [1]: import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import seaborn as sns
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
from sklearn.model_selection import train_test_split
from sklearn import preprocessing
import time
df=pd.read_csv("C:/Users/jline/Desktop/Insight/ds_1/employee_retention_data.csv")
df.head()
```

Out[1]:

	employee_id	company_id	dept	seniority	salary	join_date	quit_date
0	13021.0	7	customer_service	28	89000.0	2014-03-24	2015-10-30
1	825355.0	7	marketing	20	183000.0	2013-04-29	2014-04-04
2	927315.0	4	marketing	14	101000.0	2014-10-13	NaT
3	662910.0	7	customer_service	20	115000.0	2012-05-14	2013-06-07
4	256971.0	2	data_science	23	276000.0	2011-10-17	2014-08-22

In [2]: df.describe()

Out[2]:

	employee_id	company_id	seniority	salary
count	24702.000000	24702.000000	24702.000000	24702.000000
mean	501804.403530	3.426969	14.127803	138183.345473
std	288909.026101	2.700011	8.089520	176058.184573
min	36.000000	1.000000	1.000000	79000.000000
25%	250133.750000	1.000000	7.000000	79000.000000
50%	500793.000000	2.000000	14.000000	123000.000000
75%	753137.250000	5.000000	21.000000	187000.000000
max	999969.000000	12.000000	99.000000	408000.000000

In [3]: df.isnull().sum()

Out[3]:

	employee_id	company_id	dept	seniority	salary	join_date	quit_date
employee_id	0						
company_id	0						
dept	0						
seniority	0						
salary	0						
join_date	0						
quit_date	1132						
dtype:	int64						

In [4]: df[['join_date', 'quit_date']] = df[['join_date', 'quit_date']].apply(pd.to_datetime)

Out[4]:

	employee_id	company_id	dept	seniority	salary	join_date	quit_date
0	13021.0	7	customer_service	28	89000.0	2014-03-24	2015-10-30
1	825355.0	7	marketing	20	183000.0	2013-04-29	2014-04-04
2	927315.0	4	marketing	14	101000.0	2014-10-13	NaT
3	662910.0	7	customer_service	20	115000.0	2012-05-14	2013-06-07
4	256971.0	2	data_science	23	276000.0	2011-10-17	2014-08-22

In [5]: df['quit_bool'] = df['quit_date'].isnull()

Out[5]:

	employee_id	company_id	dept	seniority	salary	join_date	quit_date	quit_bool
0	13021.0	7	customer_service	28	89000.0	2014-03-24	2015-10-30	False
1	825355.0	7	marketing	20	183000.0	2013-04-29	2014-04-04	False
2	927315.0	4	marketing	14	101000.0	2014-10-13	NaT	True
3	662910.0	7	customer_service	20	115000.0	2012-05-14	2013-06-07	False
4	256971.0	2	data_science	23	276000.0	2011-10-17	2014-08-22	False

In [6]: a=df['quit_date']-df['join_date']

Out[6]:

	employee_id	company_id	dept	seniority	salary	join_date	quit_date	quit_bool	day
0	13021.0	7	customer_service	28	89000.0	2014-03-24	2015-10-30	False	585 days
1	825355.0	7	marketing	20	183000.0	2013-04-29	2014-04-04	False	340 days
2	927315.0	4	marketing	14	101000.0	2014-10-13	NaT	True	NaT
3	662910.0	7	customer_service	20	115000.0	2012-05-14	2013-06-07	False	389 days
4	256971.0	2	data_science	23	276000.0	2011-10-17	2014-08-22	False	1040 days

In [7]: df['day']=a

Out[7]:

	employee_id	company_id	dept	seniority	salary	join_date	quit_date	quit_bool	day
0	13021.0	7	customer_service	28	89000.0	2014-03-24	2015-10-30	False	585 days
1	825355.0	7	marketing	20	183000.0	2013-04-29	2014-04-04	False	340 days
2	927315.0	4	marketing	14	101000.0	2014-10-13	NaT	True	NaT
3	662910.0	7	customer_service	20	115000.0	2012-05-14	2013-06-07	False	389 days
4	256971.0	2	data_science	23	276000.0	2011-10-17	2014-08-22	False	1040 days

In [8]: df['time_delta_int']= df.day.astype('timedelta64[D]')

Out[8]:

	employee_id	company_id	dept	seniority	salary	join_date	quit_date	quit_bool	day	time_delta
0	13021.0	7	customer_service	28	89000.0	2014-03-24	2015-10-30	False	585 days	585.0
1	825355.0	7	marketing	20	183000.0	2013-04-29	2014-04-04	False	340 days	340.0
2	927315.0	4	marketing	14	101000.0	2014-10-13	NaT	True	NaT	NaT
3	662910.0	7	customer_service	20	115000.0	2012-05-14	2013-06-07	False	389 days	389.0
4	256971.0	2	data_science	23	276000.0	2011-10-17	2014-08-22	False	1040 days	1040.0

In [9]: join=df['join_date'].dt.month

Out[9]:

	company_id	dept	seniority	salary	quit_bool	day	time_delta_int	join_m	quit_m
0	7	customer_service	28	89000.0	False	585 days	585.0	3	10.0
1	7	marketing	20	183000.0	False	340 days	340.0	4	4.0
2	4	marketing	14	101000.0	True	NaT	NaT	10	NaT
3	7	customer_service	20	115000.0	False	389 days	389.0	5	6.0
4	2	data_science	23	276000.0	False	1040 days	1040.0	10	8.0

In [10]: #MERGE DATASET

Out[10]:

	company_id	dept	seniority	salary	quit_bool	day	time_delta_int	join_m	quit_m	quit
0	7	customer_service	28	89000.0	False	585 days	585.0	3	10.0	yes
1	7	marketing	20	183000.0	False	340 days	340.0	4	4.0	yes
2	4	marketing	14	101000.0	True	NaT	NaT	10	NaT	no
3	7	customer_service	20	115000.0	False	389 days	389.0	5	6.0	yes
4	2	data_science	23	276000.0	False	1040 days	1040.0	10	8.0	yes

In [11]: leaving=df.dropna()

Out[11]:

	dept	seniority	salary	quit_bool	day	time_delta_int	join_m	quit_m
0	customer_service	28	89000.0	False	585 days	585.0	3	10.0
1	marketing	20	183000.0	False	340 days	340.0	4	4.0
3	customer_service	20	115000.0	False	389 days	389.0	5	6.0
4	data_science	23	276000.0	False	1040 days	1040.0	10	8.0
5	data_science	14	165000.0	False	578 days	578.0	1	8.0

In [12]: sns.distplot(leaving['join_m'], hist = True, kde = True,

Out[12]:

In [13]: sns.distplot(leaving['quit_m'], hist = True, kde = True,

Out[13]:

In [14]: leaving.dtypes

Out[14]:

	dept	seniority	salary	quit_bool	day	time_delta_int	join_m	quit_m
dept	object							
seniority	int64							
salary	float64							
quit_bool	bool							
day	timedelta64[ns]							
time_delta_int	int64							
join_m	float64							
quit_m	float64							
dtype:	object							

In [15]: days=leaving['day'].dt.days

In [16]: leaving['days']= days

Out[16]:

	dept	seniority	salary	quit_bool	time_delta_int	join_m	quit_m	days	
0	customer_service	28	89000.0	False	585 days	585.0	3	10.0	585
1	marketing	20	183000.0	False	340 days	340.0	4	4.0	340
3	customer_service	20	115000.0	False	389 days	389.0	5	6.0	389
4	data_science	23	276000.0	False	1040 days	1040.0	10	8.0	1040
5	data_science	14	165000.0	False	578 days	578.0	1	8.0	578

In [16]:

Out[16]:

In [17]: sns.barplot(x = leaving['dept'].value_counts().index,

Out[17]:

In [17]: deptDummies = pd.get_dummies(df['dept'], prefix = 'dept')

Out[17]:

	dept_customer_service	dept_data_science	dept_design	dept_engineer
0	1	0	0	0
1	0	1	0	0
2	0	0	1	0
3	1	0	0	0
4	0	1	0	0
5	0	0	1	0
6	1	0	0	0
7	0	1	0	0
8	0	0	0	1
9	0	0	0	0
10	0	1	0	0
11	0	0	0	0
12	0	0	0	0
13	0	0	0	0
14	0	0	0	0
15	0	1	0	0
16	0	0	0	0
17	0	0	0	0
18	0	0	0	0
19	1	0	0	0
20	0	1	0	0
21	1	0	0	0
22	0	1	0	0
23	0	0	0	0
24	0	0	1	0
25	0	1	0	0
26	0	0	1	0
27	0	0	0	0
28	1	0	0	0
29	0	0	0	0
...
24672	0	0	0	1
24673	0	0	0	0
24674	0	0	0	0
24675	0	0	0	0
24676	0	0	0	0
24677	0	0	0	1
24678	0	0	0	0
24679	0	0	0	0
24680	0	0	0	0
24681	0	0	0	0
24682	0	0	0	1
24683	0	0	0	0
24684	0	1	0	0
24685	0	0	0	0
24686	0	0	0	0
24687	0	0	0	0
24688	0	0	0	0
24689	0	0	0	0
24690	0	0	0	0
24691	0	0	0	0
24692	0	0	1	0
24693	0	0	0	0
24694	0	0	1	0
24695	0	0	0	0
24696	0	0	0	0
24697	0	0	0	0
24698	0	0	0	0
24699	0	0	0	0
24700	0	0	0	1
24701	0	0	0	0
...
24702	rows x 6 columns			

In [18]: print(joinDummies)

Out[18]:

	joinm_0	joinm_1	joinm_2	joinm_3	joinm_4	joinm_5	joinm_6	joinm_7	joinm_8
0	0	0	0	1	0	0	0	0	0
1	0	0	0	0	1	0	0	0	0
2	0	0	0	0	0	1	0	0	0
3	0	0	0	0	0	0	1	0	0
4	0	0	0	0	0	0	0	1	0
5	1	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0
12	0	0	0	0	0	0	0	0	0
13	0	0	0	0	0	0	0	0	0
14	0	0	0	0	0	0	0	0	0
15	0	0	0	0	0	0	0	0	0
16	0	0	0	0	0	0	0	0	0
17	0	0	0	0	0	0	0	0	0
18	0	0	0	0	0	0	0	0	0
19	0	0	0	0	0	0	0	0	0
20	0	0	0	0	0	0	0	0	0
21	0	0	0	0	0	0	0	0	0
22	0	0	0	0	0	0	0	0	0
23	0	0	0	0	0	0	0	0	0
24	0	0	0	0	0	0	0	0	0
25	0	0	0	0	0	0	0	0	0
26	0	0	0	0	0	0	0	0	0
27	0	0	0	0	0	0	0	0	0
28	0	0	0	0	0	0	0	0	0
29	0	0	1	0	0	0	0	0	0
...
24672	0	0	0	0	0	0	0	0	0
24673	0	0	0	0	0				