

# IBM Human Resource Management Proposal

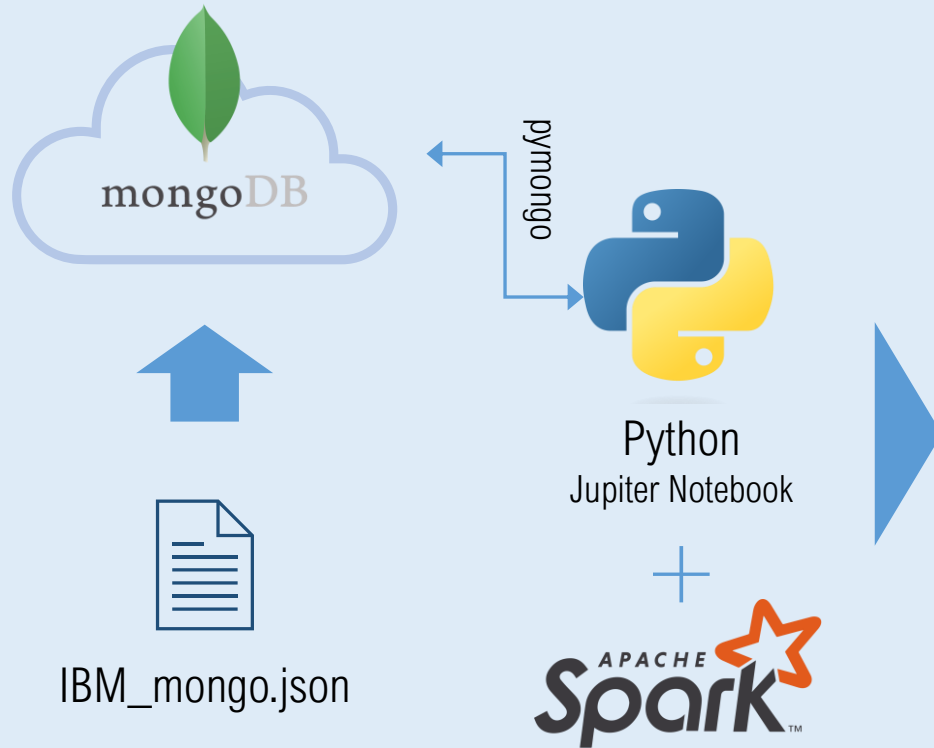




## Contents

- Data Manipulation
- Descriptive Statistics
- Hypothesis & Research Design
- Data Analysis & Visualization
- Predictive Modeling
- Conclusion

# Database Connection & Manipulation with **mongoDB**



```

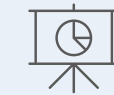
1 import pandas as pd # Pandas DataFrame
2 import seaborn as sns # Data Visualization Tool
3 import numpy as np # Numpy
4 import matplotlib.pyplot as plt # Plotting Tool
5 import statsmodels.formula.api as smf # Statistic Analysis Tool
6 from pymongo.mongo_client import MongoClient # Mongo DB Client Connection Tool
7 import urllib.parse # Url String Parsing
8
9 # setting pandas print options
10 pd.set_option('display.max_rows', 500)
11 pd.set_option('display.max_columns', 500)
12 pd.set_option('display.width', 1000)
13
14 # Connecting String
15 username = urllib.parse.quote_plus('hulyunsik')
16 password = urllib.parse.quote_plus('hulyunsik')
17 uri = "mongodb+srv://" + username + ":" + password + "@hulyunsikchoung.wfvhiwp.mongodb.net/?retryW
18
19 # Set client
20 client = MongoClient(uri)
21
22 # Send a ping to confirm a successful connection
23 try:
24     client.admin.command('ping')
25     print("Pinged your deployment. You successfully connected to MongoDB!")
26 except Exception as e:
27     print(e)

```

Pinged your deployment. You successfully connected to MongoDB!



**Data Manipulation**  
pymongo, pandas, Numpy



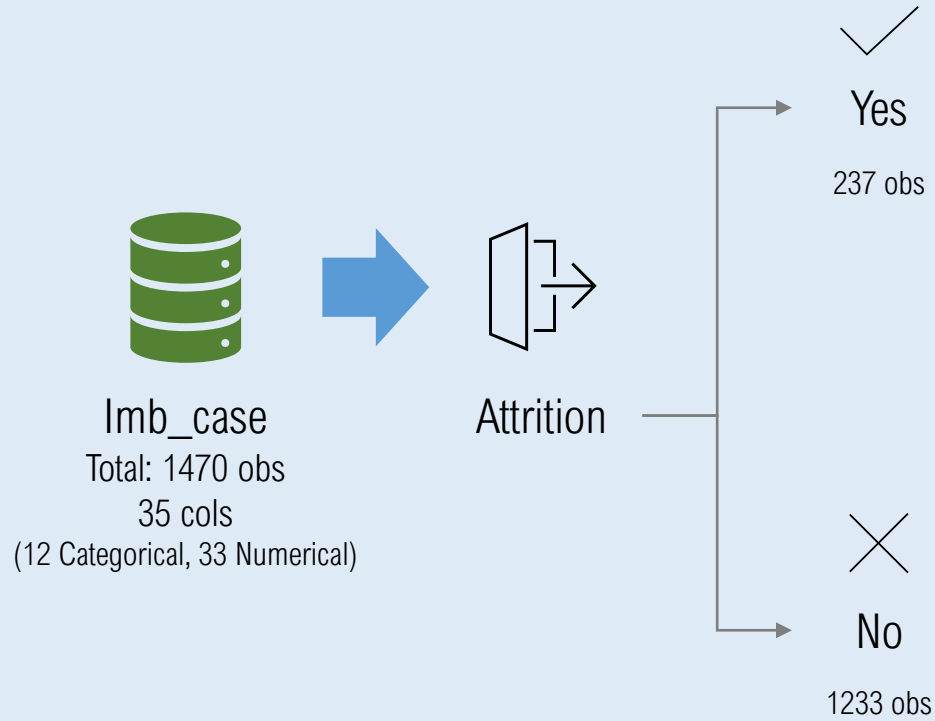
**Data Visualization**  
matplotlib, seaborn



**Data Analysis**  
Pandas, Statsmodels, PySpark

- **Data Manipulation**
- Descriptive Statistics
- Hypothesis & Research Design
- Data Analysis & Visualization
- Predictive Modeling
- Conclusion

# Descriptive Statistics Between Attrition: Numerical Variables



Attrition	Yes	No
Age	33.61	37.56
Distance From Home	10.63	8.92
Environment Satisfaction	2.46	2.77
Job Satisfaction	2.47	2.78
Monthly Income	4,787.09	6,832.74
Monthly Rate	14,559.31	14,265.78
# of Companies Worked	2.94	2.65
Percent Salary Hike	15.1	15.23
Performance Rating	3.16	3.15
Relationship Satisfaction	2.6	2.73
Stock Option Level	0.53	0.85
Total Working Years	8.24	11.86
Training Times Last Year	2.62	2.83
Work Life Balance	2.66	2.78
Years At Company	5.13	7.37
Years In Current Role	2.9	4.48
Years Since Last Promotion	1.95	2.23
Years With Current Manager	2.85	4.37

## ▼ Lower Variables

Age  
Distance from Home  
Environment Satisfaction  
Job Satisfaction  
Monthly Income  
Monthly Rate  
Relationship Satisfaction  
Stock Option Level  
Total Working Years  
Training Times Last Year  
Work Life Balance  
Years At Company  
Years In Current Role  
Years Since Last Promotion  
Years With Current Manager

## ▲ Higher Variables

# of Companies Worked

## ? Look Similar Variables

Percent Salary Hike

- Data Manipulation
- Descriptive Statistics
- Hypothesis & Research Design
- Data Analysis & Visualization
- Predictive Modeling
- Conclusion



# Descriptive Statistics Between Attrition: Categorical Variables



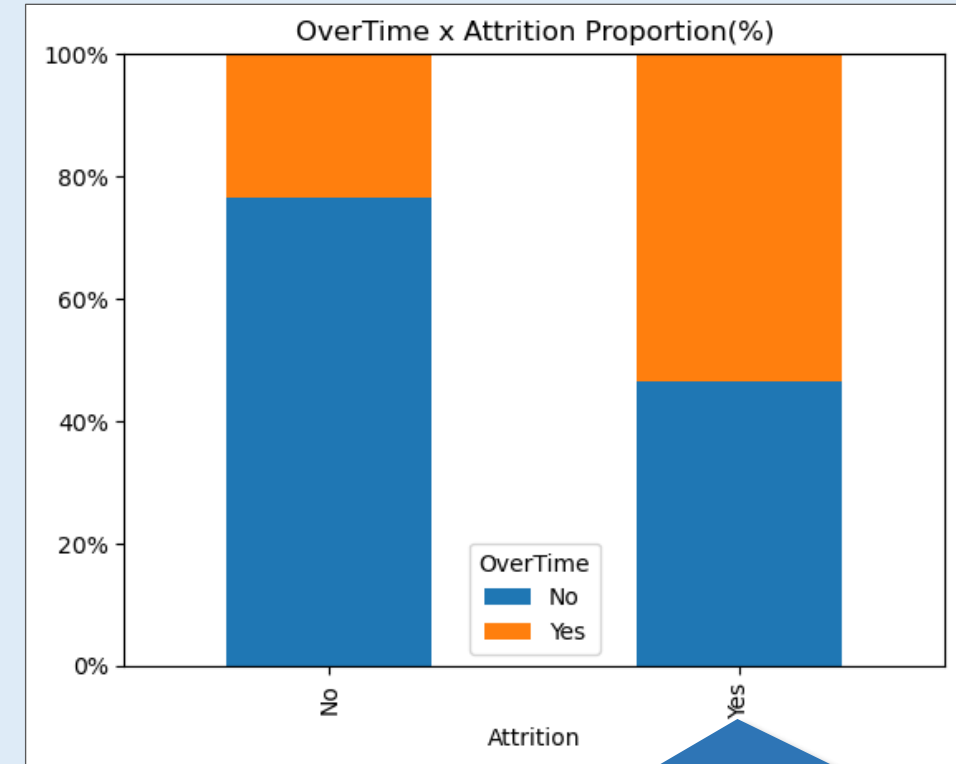
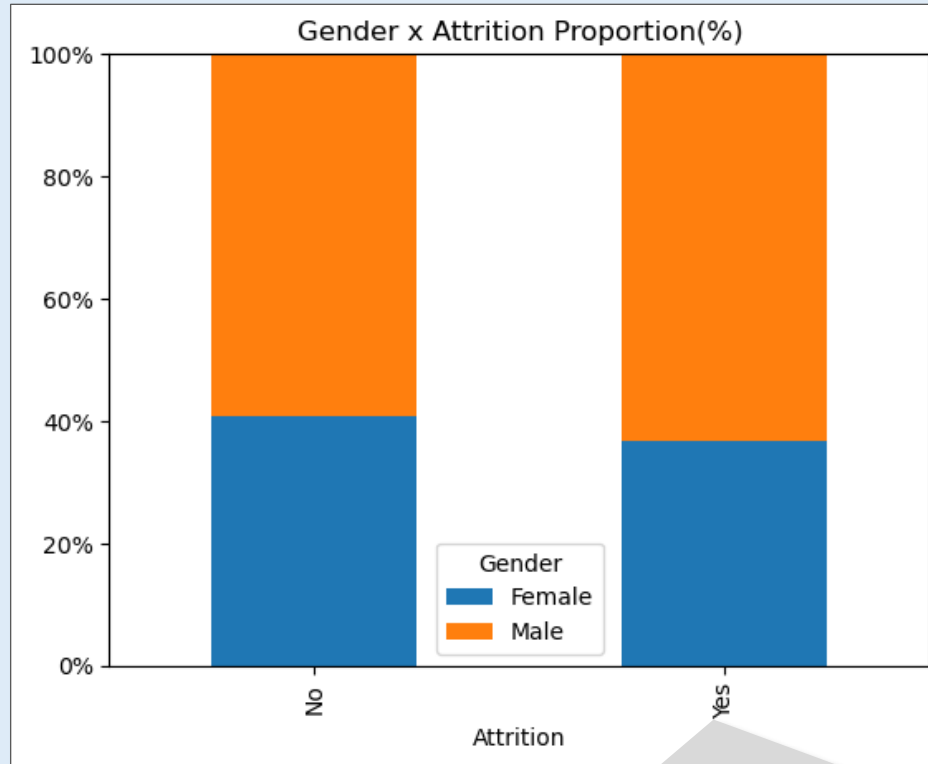
Yes

237 obs



No

1233 obs



## ▼ Less Differences within Attrition

Gender, Business Travel, Education, Job Involvement, Education Field, Over18

## ▲ More Differences within Attrition

Department, Job Role, Marital Status, Over Time

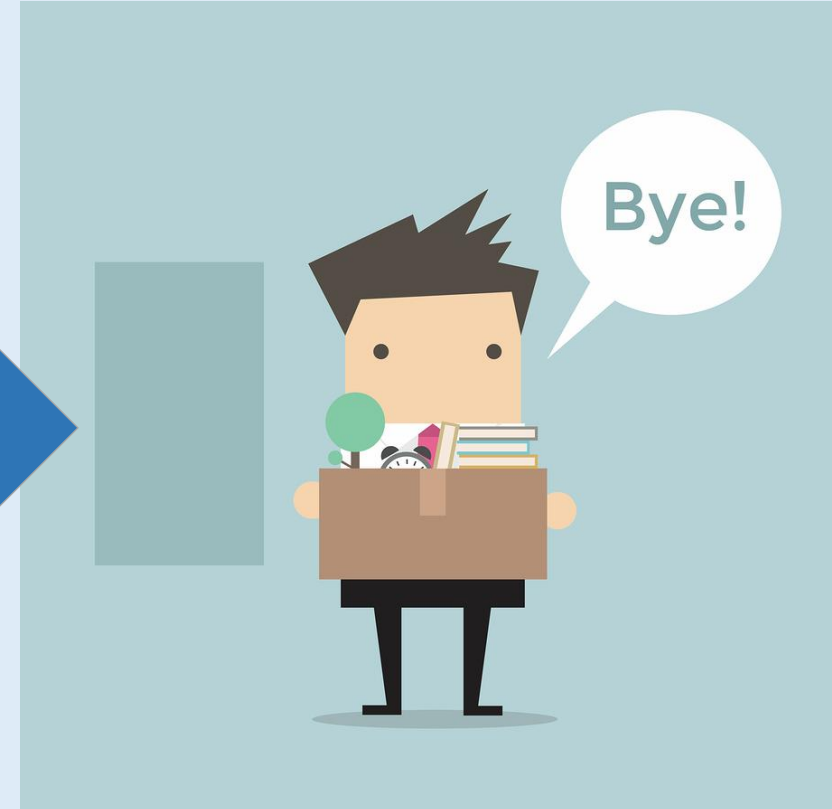
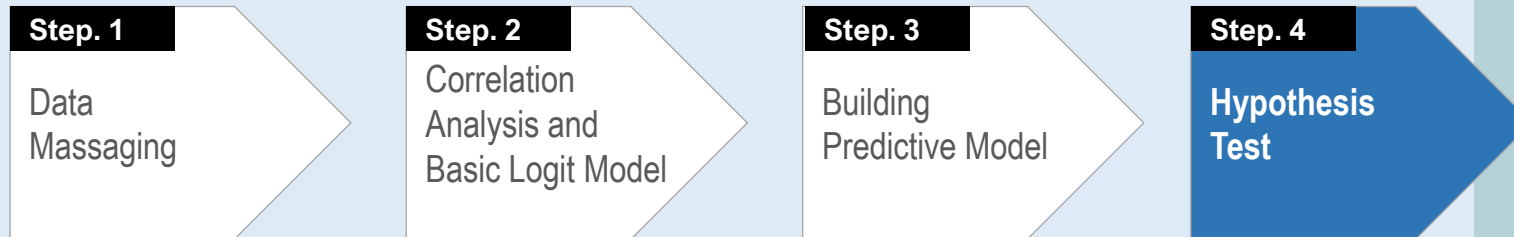
- Data Manipulation
- **Descriptive Statistics**
- Hypothesis & Research Design
- Data Analysis & Visualization
- Predictive Modeling
- Conclusion

# Hypothesis & Research Design

## Hypothesis

- $H_0$ . Variables that have more differences between Attrition status cannot predict employee attrition.  
 $H_1$ . Variables that have more differences between Attrition status can predict employee attrition.

## Research Design



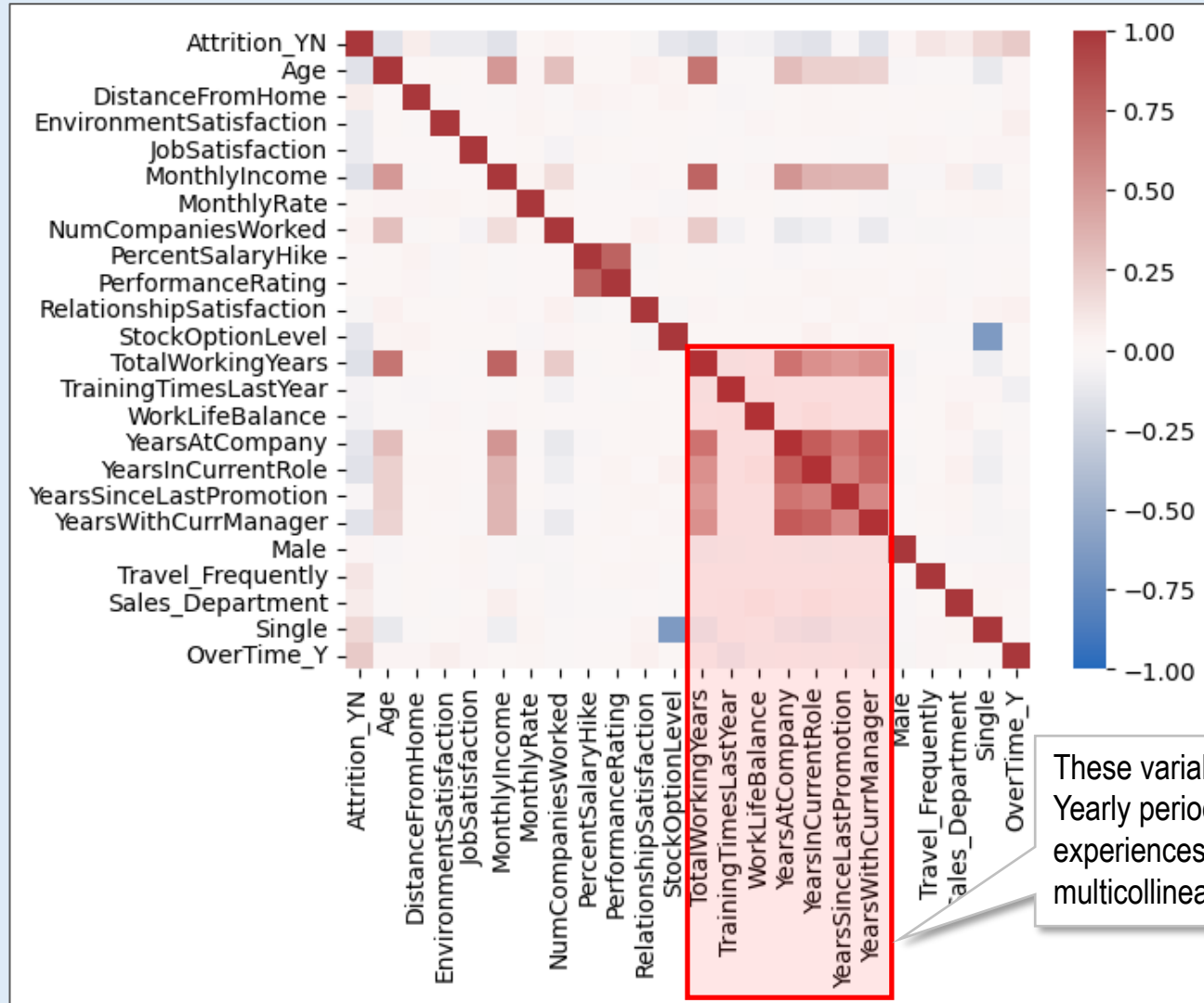
- Data Manipulation
- Descriptive Statistics
- **Hypothesis & Research Design**
- Data Analysis & Visualization
- Predictive Modeling
- Conclusion

# Data Analysis & Visualization: Correlation Analysis

Three Variables has correlation coefficient more than 0.1.

Nine Variables has negative correlation coefficient less than -0.1.

Variables	Coefficient
Over Time Yes	0.246
Single (Marital)	0.175
Travel Frequently	0.115
Sales Department	0.081
Distance From Home	0.078
Num Companies Worked	0.043
Male	0.029
Years Since Last Promotion	-0.033
Relationship Satisfaction	-0.046
Training Times Last Year	-0.059
Work Life Balance	-0.064
Environment Satisfaction	-0.103
Job Satisfaction	-0.103
Years At Company	-0.134
Stock Option Level	-0.137
Years With Current Manager	-0.156
Age	-0.159
Monthly Income	-0.160
Years In Current Role	-0.161
Total Working Years	-0.171



These variables, related in Yearly periodic working experiences, seem to have multicollinearity.

# Data Analysis & Visualization: Logistic Regression Model

## Full Model: 23 Variables

Attrition ~ Age + Distance From Home + Environment Satisfaction + Job Satisfaction + Monthly Income + Monthly Rate + Num Companies Worked + Percent Salary Hike + Performance Rating  
 + Relationship Satisfaction + Stock Option Level + Total Working Years + Training Times Last Year + Work Life Balance + Years At Company + Years In Current Role  
 + Years Since Last Promotion + Years With Current Manager + Male(Dummy) + Travel Frequently(Dummy) + Sales Department(Dummy) + Single(Dummy) + Over Time Yes(Dummy)

- Model Fitted and R-Squared is 0.281.
- Check Significance of each variables. ( $p < 0.1$ )

Model:	Logit	Pseudo R-squared:	0.281
Dependent Variable:	Attrition_YN	AIC:	981.434
#. Observations:	1470	BIC:	1,108.467
Df Model:	23	Log-Likelihood:	-466.720
Df Residuals:	1446	LL-Null:	-649.290
Converged:	1	LLR p-value:	0.000
No. Iterations:	8	Scale:	1.000

## Final Model: 17 Variables

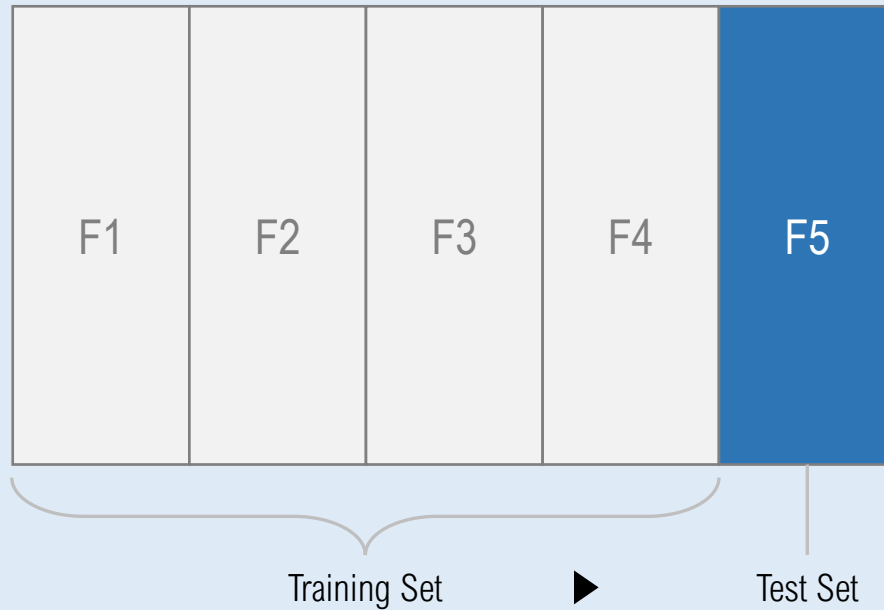
Attrition ~ Age + Distance From Home + Environment Satisfaction + Job Satisfaction + Monthly Income+ Num Companies Worked  
 + Relationship Satisfaction + Stock Option Level + Total Working Years + Training Times Last Year + Work Life Balance  
 + Years In Current Role + Male(Dummy) + Travel Frequently(Dummy) + Sales Department(Dummy) + Single(Dummy) + Over Time Yes(Dummy)



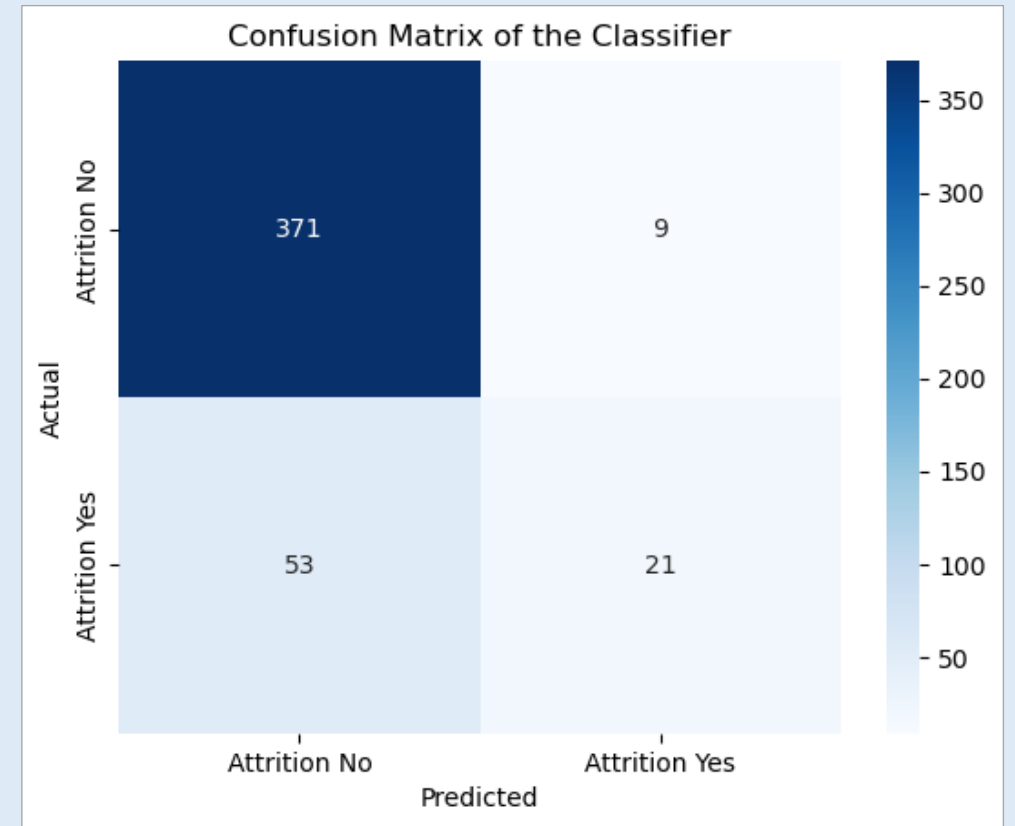
# Predictive Modeling with PySpark ML

## Cross Validating Logistic Classification ML Model

## Model Accuracy



AUC(Area Under ROC): 0.781 – Able to predict 78.1% of employees



- Data Manipulation
- Descriptive Statistics
- Hypothesis & Research Design
- Data Analysis & Visualization
- **Predictive Modeling**
- Conclusion

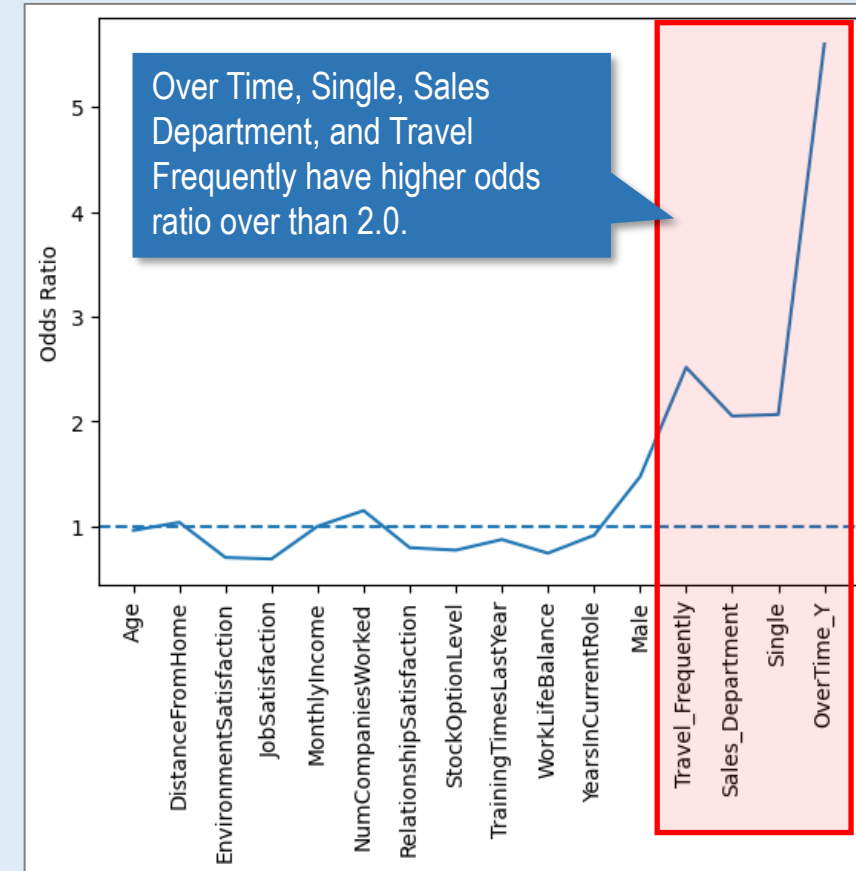
# Logistic Regression Model Result

## Final Model: 17 Variables

Attrition ~ Age + Distance From Home + Environment Satisfaction + Job Satisfaction + Monthly Income+ Num Companies Worked + Relationship Satisfaction + Stock Option Level + Total Working Years  
+ Training Times Last Year + Work Life Balance + Years In Current Role + Male(Dummy) + Travel Frequently(Dummy) + Sales Department(Dummy) + Single(Dummy) + Over Time Yes(Dummy)

Model:	Logit	Pseudo R-squared:	0.250
Dependent Variable:	Attrition_YN	AIC:	1,007.770
#. Observations:	1470	BIC:	1,097.751
Df Model:	16	Log-Likelihood:	-486.890
Df Residuals:	1453	LL-Null:	-649.290
Converged:	1	LLR p-value:	0.000
No. Iterations:	7	Scale:	1.000

Variables	Coef.	Std.Err.	P> z
Age	-0.041	0.011	0.000
DistanceFromHome	0.036	0.010	0.000
EnvironmentSatisfaction	-0.355	0.075	0.000
JobSatisfaction	-0.374	0.074	0.000
MonthlyIncome	0.000	0.000	0.000
NumCompaniesWorked	0.139	0.034	0.000
RelationshipSatisfaction	-0.229	0.076	0.003
StockOptionLevel	-0.259	0.138	0.061
TrainingTimesLastYear	-0.135	0.067	0.045
WorkLifeBalance	-0.297	0.114	0.009
YearsInCurrentRole	-0.090	0.029	0.002
Male	0.387	0.172	0.024
Travel_Frequently	0.923	0.192	0.000
Sales_Department	0.719	0.176	0.000
Single	0.726	0.226	0.001
OverTime_Y	1.723	0.174	0.000



- Data Manipulation
- Descriptive Statistics
- Hypothesis & Research Design
- Data Analysis & Visualization
- Predictive Modeling
- Conclusion

## Conclusion: How to prevent employee Attrition at IBM?



Monitoring Over Time Work



Check Working Environments of Sales Department



Single + Younger People Care



Redefine Travel Policy



Satisfaction Control – Job, Relationship, Environment Periodical Survey



IBM



Thank You



## Mongo DB Dataset

MongoDB Compass - test/Team7.imb\_case

Connect Edit View Collection Help

test Documents Team7.imb\_case

Team7.imb\_case

1.5k DOCUMENTS 1 INDEXES

Documents Aggregations Schema Indexes Validation

Filter Type a query: { field: 'value' }

EXPLAIN RESET FIND </> OPTIONS

1 - 20 of 1470

ADD DATA EXPORT DATA

SHOW 11 MORE FIELDS

```
{
  "_id": "ObjectId('64b829bfb64986d9b3edcf63')",
  "Age": 59,
  "Attrition": "No",
  "BusinessTravel": "Travel_Rarely",
  "DailyRate": 1324,
  "Department": "Research & Development",
  "DistanceFromHome": 3,
  "Education": 3,
  "EducationField": "Medical",
  "EmployeeCount": 1,
  "EmployeeNumber": 10,
  "EnvironmentSatisfaction": 3,
  "Gender": "Female",
  "HourlyRate": 81,
  "JobInvolvement": 4,
  "JobLevel": 1,
  "JobRole": "Laboratory Technician",
  "JobSatisfaction": 1,
  "MaritalStatus": "Married",
  "MonthlyIncome": 2670,
  "MonthlyRate": 9964,
  "NumCompaniesWorked": 4,
  "Over18": "Y",
  "OverTime": "Yes",
  "PercentSalaryHike": 20
}
```

```
{
  "_id": "ObjectId('64b829bfb64986d9b3edcf5d')",
  "Age": 41,
  "Attrition": "Yes",
  "BusinessTravel": "Travel_Rarely",
  "DailyRate": 1102,
  "Department": "Sales",
  "DistanceFromHome": 1,
  "Education": 2
}
```

>\_MONGOSH



## Dataset Connection with Mongo DB in local Jupyter Notebook

```
In [1]: 1 import pandas as pd # Pandas DataFrame
        2 import seaborn as sns # Data Visualization Tool
        3 import numpy as np # Numpy
        4 import matplotlib.pyplot as plt # Plotting Tool
        5 import statsmodels.formula.api as smf # Statistic Analysis Tool
        6 from pymongo.mongo_client import MongoClient # Mongo DB Client Connection Tool
        7 import urllib.parse # Url String Parsing
        8
        9 # setting pandas print options
       10 pd.set_option('display.max_rows', 500)
       11 pd.set_option('display.max_columns', 500)
       12 pd.set_option('display.width', 1000)
       13
       14 # Connecting String for MongoDB Connector
       15 username = urllib.parse.quote_plus('admin')
       16 password = urllib.parse.quote_plus('test1QW!')
       17 uri = "mongodb+srv://" + username + ":" + password + "@hultyunsikchoung.wfvhiw.mongodb.net/?retryWrites=true&w=majority"
       18
       19 # Set client
       20 client = MongoClient(uri)
       21
       22 # Send a ping to confirm a successful connection
       23 try:
       24     client.admin.command('ping')
       25     print("Pinged your deployment. You successfully connected to MongoDB!")
       26 except Exception as e:
       27     print(e)
```

Pinged your deployment. You successfully connected to MongoDB!

```
In [2]: 1 # Check Database List
        2 print(client.list_database_names())

['Team7', 'Team7_Practice', 'pokemon_Sanfran', 'sample_airbnb', 'sample_analytics', 'sample_geospatial', 'sample_guides', 'sample_mflix', 'sample_restaurants', 'sample_supplies', 'sample_training', 'sample_weatherdata', 'admin', 'local']
```

# Descriptive Statistics using Mongo DB at Local Jupyter

```
In [6]: 1 # Descriptive Statistics (Numerical Variables) using aggregate Function with MongoDB QL
2 descriptive_statistics = mycol.aggregate([
3     "$group": {
4         "_id": "$Attrition", # Attrition Status Yes / No
5         "CNT": { # Count each Attrition Status
6             "$count": {}
7         }, "Avg_Age": { # Average Age between Yes and No
8             "$avg": "$Age"
9         }, "Avg_DistanceFromHome": { # Average Distance From Home
10            "$avg": "$DistanceFromHome"
11        }, "Avg_EnvironmentSatisfaction": { # Average Environment Satisfaction Likert Scale
12            "$avg": "$EnvironmentSatisfaction"
13        }, "Avg_JobSatisfaction": { # Average Job Satisfaction Likert Scale
14            "$avg": "$JobSatisfaction"
15        }, "avg_MonthlyIncome": { # Average Monthly Income in $
16            "$avg": "$MonthlyIncome"
17        }, "Avg_Monthly_rate": { # Average Monthly Rate
18            "$avg": "$MonthlyRate"
19        }, "Avg_NumCompaniesWorked": { # Average# of Companies worked
20            "$avg": "$NumCompaniesWorked"
21        }, "Avg_PercentSalaryHike": { # Average Percent Salary hike
22            "$avg": "$PercentSalaryHike"
23        }, "Avg_RelationshipSatisfaction": { # Average Relationship Satisfaction
24            "$avg": "$RelationshipSatisfaction"
25        }, "Avg_StockOptionLevel": { # Average Stock Option Level. Looks binary 1: Yes, 0: No Stock Option
26            "$avg": "$StockOptionLevel"
27        }, "Avg_TotalWorkingYears": { # Average Total working years
28            "$avg": "$TotalWorkingYears"
29        }, "Avg_TrainingTimesLastYear": { # Average Training Times Last year
30            "$avg": "$TrainingTimesLastYear"
31        }, "Avg_WorkLifeBalance": { # Average Work Life Balance. TIME SPENT BETWEEN WORK AND OUTSIDE
32            "$avg": "$WorkLifeBalance"
33        }, "Avg_YearsAtCompany": { # Average Years at company
34            "$avg": "$YearsAtCompany"
35        }, "Avg_YearsInCurrentRole": { # Average Years in current role
36            "$avg": "$YearsInCurrentRole"
37        }, "Avg_YearsSinceLastPromotion": { # Average Years Since Last promotion
38            "$avg": "$YearsSinceLastPromotion"
39        }, "Avg_YearsWithCurrManager": { # Average Years with current manager
40            "$avg": "$YearsWithCurrManager"
41        }
42    }
43 ])
44 list(descriptive_statistics)
45 # DistanceFromHome RelationshipSatisfaction PercentSalaryHike EnvironmentSatisfaction
46 # YearsAtCompany YearsInCurrentRole YearsSinceLastPromotion YearsWithCurrManager
```

```
In [7]: 1 # Categorical Variables Crosstabulation
2 # Gender pivoting
3 descriptive_gender = mycol.aggregate([
4     "$group": {
5         "_id": ["$Attrition", "$Gender"],
6         "CNT": {
7             "$count": {}
8         }
9     }
10 })
11 list(descriptive_gender)
```

```
Out[7]: [{"_id": ['No', 'Male'], 'CNT': 732},
{'_id': ['Yes', 'Male'], 'CNT': 150},
{'_id': ['No', 'Female'], 'CNT': 501},
{'_id': ['Yes', 'Female'], 'CNT': 87}]
```

```
In [ ]: 1
```

```
In [9]: 1 # Categorical Variables Crosstabulation
2 # Department pivoting
3 descriptive_Department = mycol.aggregate([
4     "$group": {
5         "_id": ["$Attrition", "$Department"],
6         "CNT": {
7             "$count": {}
8         }
9     }
10 })
11 list(descriptive_Department)
```

```
Out[9]: [{"_id": ['Yes', 'Sales'], 'CNT': 92},
{'_id': ['No', 'Sales'], 'CNT': 354},
{'_id': ['No', 'Human Resources'], 'CNT': 51},
{'_id': ['Yes', 'Human Resources'], 'CNT': 12},
{'_id': ['No', 'Research & Development'], 'CNT': 828},
{'_id': ['Yes', 'Research & Development'], 'CNT': 13
3}]
```

```
In [8]: 1 # Categorical Variables Crosstabulation
2 # BusinessTravel pivoting
3 descriptive_BusinessTravel = mycol.aggregate([
4     "$group": {
5         "_id": ["$Attrition", "$BusinessTravel"],
6         "CNT": {
7             "$count": {}
8         }
9     }
10 })
11 list(descriptive_BusinessTravel)
```

```
Out[8]: [{"_id": ['No', 'Travel_Frequently'], 'CNT': 208},
{'_id': ['Yes', 'Travel_Rarely'], 'CNT': 156},
{'_id': ['Yes', 'Travel_Frequently'], 'CNT': 69},
{'_id': ['No', 'Non-Travel'], 'CNT': 138},
{'_id': ['Yes', 'Non-Travel'], 'CNT': 12},
{'_id': ['No', 'Travel_Rarely'], 'CNT': 887}]
```

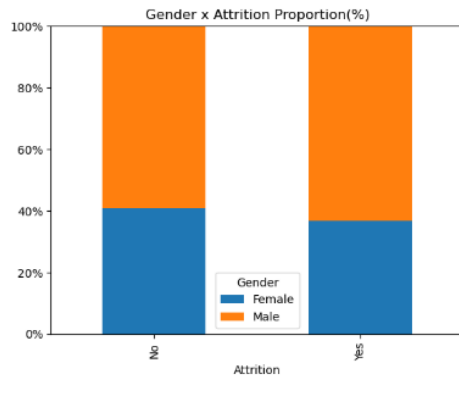
```
In [10]: 1 # Categorical Variables Crosstabulation
2 # Education pivoting
3 descriptive_Education = mycol.aggregate([
4     "$group": {
5         "_id": ["$Attrition", "$Education"],
6         "CNT": {
7             "$count": {}
8         }
9     }
10 })
11 list(descriptive_Education)
```

```
Out[10]: [{"_id": ['No', 1], 'CNT': 139},
{'_id': ['Yes', 2], 'CNT': 44},
{'_id': ['No', 2], 'CNT': 238},
{'_id': ['Yes', 3], 'CNT': 99},
{'_id': ['No', 3], 'CNT': 473},
{'_id': ['Yes', 4], 'CNT': 58},
{'_id': ['No', 5], 'CNT': 43},
{'_id': ['Yes', 1], 'CNT': 31},
{'_id': ['No', 4], 'CNT': 340},
{'_id': ['Yes', 5], 'CNT': 5}]
```

# Data Analysis & Visualization using Pandas and PyPlot

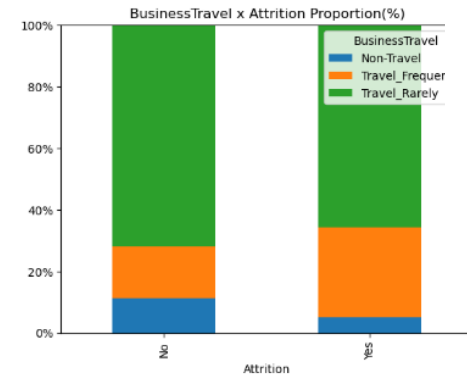
```
In [18]: 1 # Categorical Variables Crosstabulation in Pandas
2 # Visualize stacked percentage bar chart
3 col = 'Gender'
4 display(pd.crosstab(
5     columns=df['Attrition'],
6     index=df[col]
7 ))
8 pd.crosstab(
9     index=df['Attrition'],
10    columns=df[col],
11    normalize='index'
12 ).plot(kind='bar', stacked=True,
13        title=col + " x Attrition Proportion(%)")
14 plt.ylim(0,1)
15 plt.yticks(plt.yticks()[0], ['{:.0%}'.format(x) for x in plt.yticks()[0]])
16 plt.show()
```

Attrition	No	Yes
Gender		
Female	501	87
Male	732	150

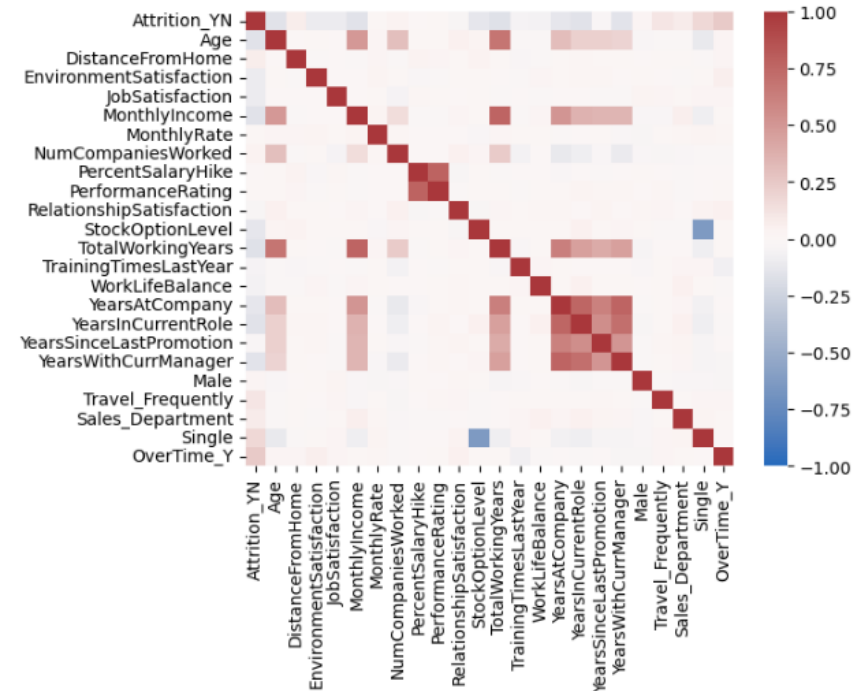


```
In [19]: 1 # Categorical Variables Crosstabulation in Pandas
2 # Visualize stacked percentage bar chart
3 col = 'BusinessTravel'
4 display(pd.crosstab(
5     columns=df['Attrition'],
6     index=df[col]
7 ))
8 pd.crosstab(
9     index=df['Attrition'],
10    columns=df[col],
11    normalize='index'
12 ).plot(kind='bar', stacked=True,
13        title=col + " x Attrition Proportion(%)")
14 plt.ylim(0,1)
15 plt.yticks(plt.yticks()[0], ['{:.0%}'.format(x) for x in plt.yticks()[0]])
16 plt.show()
```

Attrition	No	Yes
BusinessTravel		
Non-Travel	138	12
Travel_Frequently	208	69
Travel_Rarely	887	156



```
In [42]: 1 # Check Heatmap
2 sns.heatmap(data=cor_test, annot=False, vmax=1, vmin=-1, cmap='vlag')
3 plt.show()
```



# Data Massaging & Get Dummy Variables using Mongo DB

```

In [28]: 1 # Conditional Syntax: Building Dummy Variables from categorical variables
2 added_df = mycol.aggregate([
3     "$addFields": {
4         "Attrition_YN": { # Dependent Variable
5             "$switch": {
6                 "branches": [
7                     {"case": {"$eq": ["$Attrition", "Yes"]}, "then": 1},
8                     {"case": {"$eq": ["$Attrition", "No"]}, "then": 0}
9                 ], "default": 99
10            }
11        }, "Male": { # Gender Dummy Variable: Base = Female
12            "$switch": {
13                "branches": [
14                    {"case": {"$eq": ["$Gender", "Male"]}, "then": 1},
15                ], "default": 0
16            }
17        }, "Travel_Frequently": { # BusinessTravel Dummy Variable. Only focused on Travel Frequently
18            "$switch": {
19                "branches": [
20                    {"case": {"$eq": ["$BusinessTravel", "Travel_Frequently"]}, "then": 1}
21                ], "default": 0
22            }
23        }, "Sales_Department": { # Department Dummy Variable. Only Sales_Department
24            "$switch": {
25                "branches": [
26                    {"case": {"$eq": ["$Department", "Sales"]}, "then": 1}
27                ], "default": 0
28            }
29        }, "Single": { # MaritalStatus Dummy Variable. Only single
30            "$switch": {
31                "branches": [
32                    {"case": {"$eq": ["$MaritalStatus", "Single"]}, "then": 1}
33                ], "default": 0
34            }
35        }, "OverTime_Y": { # OverTime Binary Dummy. Over Time Yes.
36            "$switch": {
37                "branches": [
38                    {"case": {"$eq": ["$OverTime", "Yes"]}, "then": 1}
39                ], "default": 0
40            }
41        }
42    }
43 })
  
```

# Logistic Regression with Statsmodels Library at local Jupyter

```
In [47]: 1 model2 = smf.logit(data=df, formula=fml2)
2 res = model2.fit()
3
4 # Final Model
5 res.summary2()
```

Optimization terminated successfully.  
Current function value: 0.331214  
Iterations 7

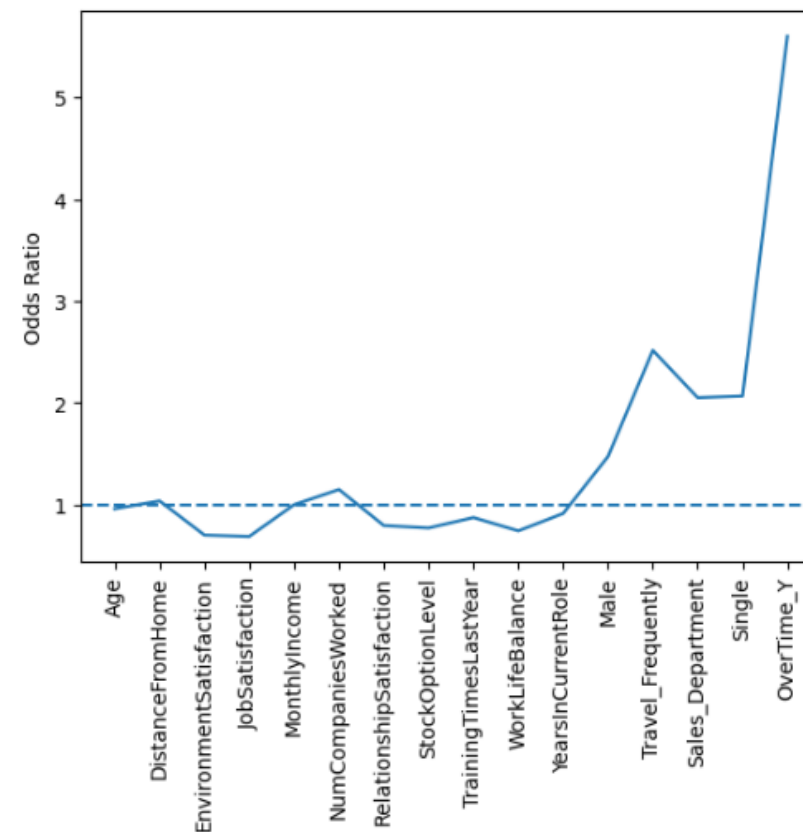
Out[47]:

Model:	Logit	Pseudo R-squared:	0.250
Dependent Variable:	Attrition_YN	AIC:	1007.7700
Date:	2023-07-25 09:51	BIC:	1097.7513
No. Observations:	1470	Log-Likelihood:	-486.89
Df Model:	16	LL-Null:	-649.29
Df Residuals:	1453	LLR p-value:	1.8140e-59
Converged:	1.0000	Scale:	1.0000
No. Iterations:	7.0000		

	Coef.	Std.Err.	z	P> z	[0.025	0.975]
Intercept	2.2186	0.6452	3.4388	0.0006	0.9541	3.4831
Age	-0.0409	0.0113	-3.6310	0.0003	-0.0630	-0.0188
DistanceFromHome	0.0363	0.0099	3.6564	0.0003	0.0168	0.0557
EnvironmentSatisfaction	-0.3545	0.0754	-4.7020	0.0000	-0.5023	-0.2067
JobSatisfaction	-0.3740	0.0744	-5.0284	0.0000	-0.5197	-0.2282
MonthlyIncome	-0.0001	0.0000	-3.6660	0.0002	-0.0002	-0.0000
NumCompaniesWorked	0.1391	0.0339	4.1027	0.0000	0.0726	0.2055
RelationshipSatisfaction	-0.2290	0.0761	-3.0104	0.0026	-0.3781	-0.0799
StockOptionLevel	-0.2593	0.1381	-1.8770	0.0605	-0.5300	0.0115
TrainingTimesLastYear	-0.1352	0.0673	-2.0076	0.0447	-0.2671	-0.0032
WorkLifeBalance	-0.2972	0.1137	-2.6136	0.0090	-0.5201	-0.0743
YearsInCurrentRole	-0.0901	0.0293	-3.0795	0.0021	-0.1474	-0.0327
Male	0.3872	0.1715	2.2582	0.0239	0.0511	0.7233
Travel_Frequently	0.9226	0.1918	4.8103	0.0000	0.5467	1.2986
Sales_Department	0.7185	0.1760	4.0813	0.0000	0.3734	1.0635
Single	0.7258	0.2263	3.2076	0.0013	0.2823	1.1694
OverTime_Y	1.7231	0.1742	9.8941	0.0000	1.3817	2.0644

```
In [50]: 1 # Odds Ratio Visualization
2 g = sns.lineplot(data=or_tbl.iloc[1:], ['Odds Ratio'], )
3 g.tick_params(axis='x', rotation=90)
4 g.axhline(y=1.0, linestyle='--')
5 plt.show()
```





# Data Store at Dataiku DFS

yunsik\_choung

Datasets

Design Node

Search DSS...

☐


+ NEW DATASET ▾

Last modified ▾

Tags ▾

Favorites

6 Datasets

☐

IBM\_final\_df\_dfs

Amazon S3 | Modified 1 minute ago

☐

IBM\_final\_df

Uploaded Files | Modified 1 minute ago

☐

Pokemon\_yunsik\_Choung\_dist

Amazon S3 | Modified 18 hours ago

☐

Pokemon\_yunsik\_Choung\_dfs

Uploaded Files | Modified 18 hours ago

☐

pokemon\_battles\_yunsik\_choung

MySQL | Modified 12 days ago

☐

pokemon\_yunsik\_choung

MySQL | Modified 12 days ago

→

Click an item to view details

# Dataiku Notebook Dataset Connection

IBM\_ychoung@student.hult.edu's Python notebook

{ } CODE SAMPLES + CREATE RECIPE C FORCE RELOAD ACTIONS

jupyter IBM\_ychoung@student.hult.edu's Python notebook Last Checkpoint: 8 minutes ago (unsaved changes)

File Edit View Insert Cell Kernel Widgets Help Trusted Python 3

In [1]:

1 %pylab inline

Populating the interactive namespace from numpy and matplotlib

In [2]:

1 import dataiku  
2 import dataiku.spark as dkuspark  
3 import pyspark  
4 from pyspark.sql import SQLContext

In [3]:

1 # Load PySpark  
2 sc = pyspark.SparkContext.getOrCreate()  
3 sqlContext = SQLContext(sc)

/home/dataiku/dataiku-dss-12.0.0/spark-standalone-home/python/pyspark/sql/context.py:114: FutureWarning: Deprecated in 3.0.0. Use SparkSession.builder.getOrCreate() instead.  
FutureWarning,

In [4]:

1 # Example: Read the descriptor of a Dataiku dataset  
2 mydataset = dataiku.Dataset("IBM\_final\_df\_dfs")  
3 # And read it as a Spark dataframe  
4 df = dkuspark.get\_dataframe(sqlContext, mydataset)

/home/dataiku/dataiku-dss-12.0.0/spark-standalone-home/python/pyspark/sql/dataframe.py:127: UserWarning: DataFrame constructor or is internal. Do not directly use it.  
warnings.warn("DataFrame constructor is internal. Do not directly use it.")

In [5]:

1 # Example: Get the count of records in the dataframe  
2 df.count()

Out[5]: 1470

# Building Logistic Classification ML Model at Dataiku PySpark

```
In [6]: 1 # Loading Libraries
2 from pyspark.ml import Pipeline
3 from pyspark.ml.classification import LogisticRegression
4 from pyspark.ml.feature import HashingTF, Tokenizer
5 from pyspark.sql import Row
6 from pyspark.sql.functions import UserDefinedFunction
7 from pyspark.ml.linalg import Vectors
8 from pyspark.ml.evaluation import BinaryClassificationEvaluator
9 from pyspark.ml.metrics import AreaUnderROC
10 from pyspark.ml.metrics import AreaUnderPRC
```

```
In [14]: 1 # Define features
2 X = ['Age', 'Distance', 'Environment', 'JobSatisfaction', 'MonthlyIncome', 'NumCompaniesWorked', 'Relationship', 'StockOptions', 'TrainingTimeInHoursPerWeek', 'WorkLifeBalance', 'YearsSinceLastJobChange', 'Male', 'TravelFrequent', 'SalesDepartment', 'Single', 'OverTime']

In [15]: 1 #creating vectors with names of variables
2 vecAssembler = VectorAssembler(inputCols = X, outputCol = 'features')
3 v_df = vecAssembler.transform(df)
4 vhouse_df = v_df.select(['features', 'Attrition_YN'])
5 vhouse_df = vhouse_df.withColumnRenamed("Attrition_YN", "label")

In [16]: 1 #splitting the dataset
2 splits = vhouse_df.randomSplit([0.7, 0.3])
3 train_df = splits[0]
4 test_df = splits[1]

In [17]: 1 #creating an object with the Logistic regression engine
2 lr = LogisticRegression(maxIter=20)
3 pipeline = Pipeline(stages=[lr])
4
5 #fitting the model
6 model = lr.fit(train_df)
7
8 #evaluating the model using testing data
9 result = model.transform(test_df)
10 result.show()
11
12
13 from pyspark.ml.evaluation import BinaryClassificationEvaluator
14 evaluator = BinaryClassificationEvaluator(rawPredictionCol="rawPrediction")
15 AUC_ROC = evaluator.evaluate(result, {evaluator.metricName: "areaUnderROC"})
16 print('AUC ROC:' + str(AUC_ROC))
```

features	label	rawPrediction	probability	prediction
(16,[0,1,2,3,4,5,...])	0	[2.38577825943823...	[0.91573637605174...	0.0
(16,[0,1,2,3,4,5,...])	0	[3.09257914366593...	[0.95658560204141...	0.0
(16,[0,1,2,3,4,5,...])	0	[4.51570331595847...	[0.98918238940619...	0.0
(16,[0,1,2,3,4,5,...])	1	[2.93904835883414...	[0.94974332351530...	0.0
(16,[0,1,2,3,4,5,...])	1	[6.80194334287698...	[0.99888962150833...	0.0
(16,[0,1,2,3,4,6,...])	0	[3.21327987735179...	[0.96133097547455...	0.0
(16,[0,1,2,3,4,6,...])	0	[1.69647147146948...	[0.84507332672845...	0.0
(16,[0,1,2,3,4,6,...])	1	[2.29889215270991...	[0.90878524598189...	0.0
(16,[0,1,2,3,4,6,...])	0	[2.88163541994249...	[0.94693110777703...	0.0
(18,0,5,0,2,0,2,0...	1	[-2.6484791853138...	[0.06608280597565...	1.0
(18,0,5,0,2,0,4,0...	0	[1.32389054697974...	[0.78982826391330...	0.0
(18,0,10,0,4,0,3,...)	0	[0.40444001655082...	[0.59975395282569...	0.0
(19,0,21,0,4,0,2,...)	1	[-1.6146692393790...	[0.16594135990624...	1.0
(20,0,2,0,3,0,3,0...	1	[1.39588010226257...	[0.80152930900454...	0.0
(20,0,2,0,3,0,3,0...	0	[1.05199964666936...	[0.74115870131374...	0.0
(20,0,3,0,1,0,3,0...	0	[0.04981006246190...	[0.51244994165256...	0.0
(20,0,4,0,1,0,1,0...	1	[-0.1686360448166...	[0.45794061572145...	1.0
(20,0,9,0,4,0,1,0...	0	[-0.1409727814251...	[0.46481505550586...	1.0
(20,0,10,0,4,0,3,...)	1	[0.19399076441069...	[0.54834617084121...	0.0
(21,0,15,0,3,0,4,...)	0	[1.88839258608584...	[0.86857214559722...	0.0

only showing top 20 rows

AUC ROC:0.8098939357137055

# Cross-Validating Logistic Classification ML Model at Dataiku PySpark

```
In [18]: 1 from pyspark.ml.evaluation import BinaryClassificationEvaluator
2
3 # Evaluate model
4 evaluator = BinaryClassificationEvaluator(rawPredictionCol="rawPrediction")
5 evaluator.evaluate(result)
6
7 from pyspark.ml.tuning import ParamGridBuilder, CrossValidator
8
9 # Create ParamGrid for Cross Validation
10 paramGrid = (ParamGridBuilder()
11             .addGrid(lr.regParam, [0.01, 0.5, 2.0])
12             .addGrid(lr.elasticNetParam, [0.0, 0.5, 1.0])
13             .addGrid(lr.maxIter, [1, 5, 10])
14             .build())
15
16 # Create 5-fold CrossValidator
17 cv = CrossValidator(estimator=lr, estimatorParamMaps=paramGrid, evaluator=evaluator, numFolds=5)
18
19 # Run cross validations
20 cvModel = cv.fit(train_df)
21 # this will likely take a fair amount of time because of the amount of models that we're training
22
23 # Use test set to measure the accuracy of our model on new data
24 predictions = cvModel.transform(test_df)
25
26 # cvModel uses the best model found from the Cross Validation
27 # Evaluate best model
28 evaluator.evaluate(predictions)
29
30 print('Model Intercept: ', cvModel.bestModel.intercept)
31 weights = cvModel.bestModel.coefficients
32 weights = [(float(w),) for w in weights] # convert numpy type to float, and to tuple
33 weightsDF = sqlContext.createDataFrame(weights, ["Feature Weight"])
34 weightsDF.show()
35 # View best model's predictions and probabilities of each prediction class
36 selected = predictions.select("label", "prediction", "probability", "features")
37 selected.show()
```

Model Intercept: 1.9190516483834124

Feature Weight	label	prediction	probability	features
-0.04369386980471...	0	0.0	[0.89928616669743...	(16,[0,1,2,3,4,5,...
0.035427510733245716	0	0.0	[0.94272640413864...	(16,[0,1,2,3,4,5,...
-0.2551387647715912	0	0.0	[0.98182615401102...	(16,[0,1,2,3,4,5,...
-0.27628731814289725	1	0.0	[0.93367768167388...	(16,[0,1,2,3,4,5,...
-9.06374172468278...	1	0.0	[0.99737433867805...	(16,[0,1,2,3,4,5,...
0.12822677528119478	0	0.0	[0.94581082995688...	(16,[0,1,2,3,4,6,...
-0.23009935781559307	0	0.0	[0.83367294845035...	(16,[0,1,2,3,4,6,...
-0.14807598935330474	1	0.0	[0.89169093077322...	(16,[0,1,2,3,4,6,...
-0.10791676899760719	0	0.0	[0.93043306011001...	(16,[0,1,2,3,4,6,...
-0.35650727402803806	1	1.0	[0.10398648737237...	[18.0,5.0,2.0,2.0...
-0.10134574674295341	0	0.0	[0.78073104674648...	[18.0,5.0,2.0,4.0...
0.4458859863296991	0	0.0	[0.62089237210116...	[18.0,10.0,4.0,3...
0.9939280604547683	1	1.0	[0.22195636616613...	[19.0,21.0,4.0,2...
0.6025395829110497	1	0.0	[0.79082470724875...	[20.0,2.0,3.0,3.0...
0.7626955650098408	0	0.0	[0.74328136181085...	[20.0,2.0,3.0,3.0...
1.45745781253733	0	0.0	[0.55780977389199...	[20.0,3.0,1.0,3.0...
	1	0.0	[0.50568789919049...	[20.0,4.0,1.0,1.0...
	0	0.0	[0.51669241271838...	[20.0,9.0,4.0,1.0...
	1	0.0	[0.56872664126864...	[20.0,10.0,4.0,3...
	0	0.0	[0.85194641494331...	[21.0,15.0,3.0,4...

only showing top 20 rows