simplilearn Assignment: IBM HR Analytics Employee Attrition Modeling The comments/sections provided are your cues to perform the assignment. You don't need to limit yourself to the number of rows/cells provided. You can add additional rows in each section to add more lines of code. If at any point in time you need help on solving this assignment, view our demo video to understand the different steps of the code. Happy coding! IBM HR Analytics Employee Attrition Modeling. **DESCRIPTION** IBM is an American MNC operating in around 170 countries with major business vertical as computing, software, and hardware. Attrition is a major risk to service-providing organizations where trained and experienced people are the assets of the company. The organization would like to identify the factors which influence the attrition of employees. **Data Dictionary:** Age: Age of employee Attrition: Employee attrition status Department: Department of work DistanceFromHome Education: 1-Below College; 2- College; 3-Bachelor; 4-Master; 5-Doctor; EducationField EnvironmentSatisfaction: 1-Low; 2-Medium; 3-High; 4-Very High; JobSatisfaction: 1-Low; 2-Medium; 3-High; 4-Very High; MaritalStatus MonthlyIncome NumCompaniesWorked: Number of companies worked prior to IBM WorkLifeBalance: 1-Bad; 2-Good; 3-Better; 4-Best; YearsAtCompany: Current years of service in IBM **Analysis Task:** • Import attrition dataset and import libraries such as pandas, matplotlib.pyplot, numpy, and seaborn. • Exploratory data analysis Find the age distribution of employees in IBM Explore attrition by age Explore data for Left employees Find out the distribution of employees by the education field Give a bar chart for the number of married and unmarried employees Build up a logistic regression model to predict which employees are likely to attrite. **Table of Contents** 1 Import libraries and dataset 2 Exploratory data analysis 2.1 Age distribution of Employees 2.2 Explore attrition by age 2.3 Explore data for Left employees 2.4 Find out the distribution of employees by the education field 2.5 Give a bar chart for the number of married and unmarried employees • 3 Build up a logistic regression model to predict which employees are likely to attrite 3.1 Data Preprocessing 3.2 Train and Test data 3.3 Build and Evaluate Model Import libraries and dataset import numpy as np import pandas as pd import matplotlib.pyplot as plt %matplotlib inline from patsy import dmatrices import sklearn import seaborn as sns dataframe=pd.read csv("D:\\NIPUN SC REC\\3 Practice Project\\Course 5 Data Science with **Exploratory data analysis** dataframe.head() DistanceFromHome **Education EducationField EnvironmentSatisfaction** JobSa Age Attrition Department 41 2 Life Sciences 2 0 Yes Sales 1 Research & 49 Life Sciences No Development Research & 2 37 2 2 Other Yes Development Research & 33 Life Sciences 3 No Development Research & 27 2 1 Medical 1 No Development In [4]: names = dataframe.columns.values print(names) ['Age' 'Attrition' 'Department' 'DistanceFromHome' 'Education' 'EducationField' 'EnvironmentSatisfaction' 'JobSatisfaction' 'MaritalStatus' 'MonthlyIncome' 'NumCompaniesWorked' 'WorkLifeBalance' 'YearsAtCompany'] dataframe.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 1470 entries, 0 to 1469 Data columns (total 13 columns): Column Non-Null Count Dtype 1470 non-null 0 Age int64 Attrition 1470 non-null object 1470 non-null object Department DistanceFromHome 1470 non-null int64 Education 1470 non-null int64 EducationField 1470 non-null object EnvironmentSatisfaction 1470 non-null int64 JobSatisfaction 1470 non-null int64 MaritalStatus 1470 non-null objec MaritalStatus 1470 non-null object 1470 non-null int64 MonthlyIncome 10 NumCompaniesWorked 1470 non-null int64 11 WorkLifeBalance 1470 non-null int64 12 YearsAtCompany 1470 non-null int64 dtypes: int64(9), object(4) memory usage: 126.4+ KB dataframe.shape Out[6]: (1470, 13) The Attrition dataset has 1470 observations with 13 variables. Out of the 13 variables, there exists one target variable 'Attrition' with possible outcomes Yes and No. The other 12 variables are independent variables. Age distribution of Employees # histogram for age plt.figure(figsize=(7,5)) dataframe['Age'].hist(bins=70) plt.title("Age distribution of Employees") plt.xlabel("Age") plt.ylabel("# of Employees") plt.show() Age distribution of Employees 80 60 of Employees 50 40 30 20 10 Explore attrition by age #Explore data for Attrition by Age plt.figure(figsize=(14,10)) plt.scatter(dataframe.Attrition,dataframe.Age, alpha=0.1) plt.title("Attrition by Age ") plt.ylabel("Age") plt.grid(b=True, which='major',axis='y') plt.show() Attrition by Age 50 30 20 The Scatter plot shows most of the attritions are centered around 30 year of age. **Explore data for Left employees** dataframe.Attrition.value counts() 1233 No 237 Name: Attrition, dtype: int64 The dataset is well organised with no missing values. This is a Binary Classification Problem, so the Distribution of instances among the 2 classes is visualized below. # explore data for Left employees breakdown plt.figure(figsize=(8,6)) dataframe.Attrition.value counts().plot(kind='bar',color='blue',alpha=.65) plt.title("Attrition breakdown") plt.show() Attrition breakdown 1200 1000 800 600 400 200 lés S The above plot shows the distribution of Attrition. Out of the total of 1470 observations 1233 is No, whereas 237 is Yes. We will treat this imbalance after splitting the data into Training and Test Set. Find out the distribution of employees by the education field # explore data for Education Field distribution plt.figure(figsize=(7,5)) dataframe.EducationField.value_counts().plot(kind='barh',color='g',alpha=.65) plt.title("Education Field Distribution") plt.show() Education Field Distribution Human Resources Other Technical Degree Marketing Medical Life Sciences Ò 100 200 300 400 500 600 **Attrition by Education Field** pd.crosstab(dataframe["EducationField"],dataframe["Attrition"]).plot(kind="barh",stacl plt.title("Attrition by Education Field") plt.xlabel("Education Field") plt.ylabel("Frequency of Attrition") plt.show() Attrition by Education Field Attrition Technical Degree Νo Yes Other Frequency of Attrition Medical Marketing Life Sciences Human Resources 100 200 300 400 500 600 Education Field LifeSciences and Medical fields have the highest number of employees and highest number of attrition rate. The percentage of employees who have attritioned against those who have been retained seems to be approximately same in all the education fields. Give a bar chart for the number of married and unmarried employees # explore data for Marital Status plt.figure(figsize=(7,5)) dataframe.MaritalStatus.value counts().plot(kind='bar',alpha=.5) plt.show() 700 600 500 400 300 200 100 0 **Attrition by Marital Status** In [14]: pd.crosstab(dataframe["MaritalStatus"],dataframe["Attrition"]).plot(kind="bar",stacked plt.title("Attrition by Marital Status") plt.xlabel("Marital Status") plt.ylabel("Frequency of Attrition") plt.show() Attrition by Marital Status 600 Attrition No 500 Yes Frequency of Attrition 400 300 200 100 0 Married Marital Status The highest attrition is well correlated to 'Single' followed by 'Married' & 'Divorced. Build up a logistic regression model to predict which employees are likely to attrite dataframe.describe() DistanceFromHome Education **EnvironmentSatisfaction JobSatisfaction** MonthlyIncome 1470.000000 count 1470.000000 1470.000000 1470.000000 1470.000000 1470.000000 mean 36.923810 9.192517 2.912925 2.721769 2.728571 6502.931293 std 9.135373 8.106864 1.024165 1.093082 1.102846 4707.956783 min 18.000000 1.000000 1.000000 1.000000 1.000000 1009.000000 30.000000 25% 2.000000 2.000000 2.000000 2.000000 2911.000000 **50%** 36.000000 7.000000 3.000000 3.000000 3.000000 4919.000000 **75%** 43.000000 14.000000 4.000000 4.000000 4.000000 8379.000000 60.000000 29.000000 5.000000 4.000000 4.000000 19999.000000 dataframe.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 1470 entries, 0 to 1469 Data columns (total 13 columns): Non-Null Count Dtype Age 1470 non-null int64 0 1470 non-null object Attrition 1470 non-null object 1470 non-null int64 Department DistanceFromHome 1470 non-null int64 Education EducationField 1470 non-null object EnvironmentSatisfaction 1470 non-null int64 JobSatisfaction 1470 non-null int64 MaritalStatus 1470 non-null object 8 1470 non-null int64 MonthlyIncome 1470 non-null 10 NumCompaniesWorked int64 11 WorkLifeBalance 1470 non-null int64 1470 non-null int64 12 YearsAtCompany dtypes: int64(9), object(4) memory usage: 126.4+ KB dataframe.columns Out[17]: Index(['Age', 'Attrition', 'Department', 'DistanceFromHome', 'Education', 'EducationField', 'EnvironmentSatisfaction', 'JobSatisfaction', 'MaritalStatus', 'MonthlyIncome', 'NumCompaniesWorked', 'WorkLifeBalance', 'YearsAtCompany'], dtype='object') dataframe.std() 9.135373 Out[18]: Age 8.106864 DistanceFromHome Education 1.024165 EnvironmentSatisfaction 1.093082 JobSatisfaction 1.102846 4707.956783 MonthlyIncome NumCompaniesWorked 2.498009 WorkLifeBalance 0.706476 6.126525 YearsAtCompany dtype: float64 dataframe['Attrition'].value counts() Out[19]: No 1233 Yes 237 Name: Attrition, dtype: int64 dataframe['Attrition'].dtypes Out[20]: dtype('0') Data Preprocessing dataframe['Attrition'].replace('Yes',1, inplace=True) dataframe['Attrition'].replace('No',0, inplace=True) dataframe.head(10) Age Attrition Department DistanceFromHome Education EducationField EnvironmentSatisfaction JobSa 41 Sales Life Sciences Research & 49 Life Sciences 1 Development Research & 2 37 Other Development Research & Life Sciences 3 33 Development Research & 4 27 2 1 Medical 1 Development Research & Life Sciences 5 32 Development Research & 3 3 6 59 3 Medical Development Research & Life Sciences 30 7 24 1 Development Research & 8 38 23 3 Life Sciences Development Research & 9 36 27 Medical Development # building up a logistic regression model X = dataframe.drop(['Attrition'],axis=1) X.head() Y = dataframe['Attrition'] Y.head() 0 1 0 1 3 0 Name: Attrition, dtype: int64 In [24]: dataframe['EducationField'].value_counts() Out[24]: Life Sciences 606 Medical 464 Marketing 159 Technical Degree 132 Other Human Resources 27 Name: EducationField, dtype: int64 dataframe['EducationField'].replace('Life Sciences',1, inplace=True) dataframe['EducationField'].replace('Medical',2, inplace=True) dataframe['EducationField'].replace('Marketing', 3, inplace=True) dataframe['EducationField'].replace('Other',4, inplace=True) dataframe['EducationField'].replace('Technical Degree',5, inplace=True) dataframe['EducationField'].replace('Human Resources', 6, inplace=True) dataframe['EducationField'].value counts() 606 464 3 159 5 132 4 82 Name: EducationField, dtype: int64 dataframe['Department'].value counts() Research & Development Sales 446 Human Resources Name: Department, dtype: int64 dataframe['Department'].replace('Research & Development',1, inplace=True) dataframe['Department'].replace('Sales',2, inplace=True) dataframe['Department'].replace('Human Resources', 3, inplace=True) dataframe['Department'].value_counts() 961 446 63 Name: Department, dtype: int64 dataframe['MaritalStatus'].value counts() 673 Out[30]: Married Single 470 Divorced 327 Name: MaritalStatus, dtype: int64 dataframe['MaritalStatus'].replace('Married',1, inplace=True) dataframe['MaritalStatus'].replace('Single',2, inplace=True) dataframe['MaritalStatus'].replace('Divorced',3, inplace=True) dataframe['MaritalStatus'].value counts() Out[32]: 1 470 327 Name: MaritalStatus, dtype: int64 x=dataframe.select dtypes(include=['int64']) x.dtypes Out[33]: Age int64 Attrition int64 int64 Department DistanceFromHome int64 Education int64 EducationField int64 EnvironmentSatisfaction int64 JobSatisfaction int64 MaritalStatus int64 int64 MonthlyIncome int64 NumCompaniesWorked int64 WorkLifeBalance YearsAtCompany int64 dtype: object In [34]: x.columns Out[34]: Index(['Age', 'Attrition', 'Department', 'DistanceFromHome', 'Education', 'EducationField', 'EnvironmentSatisfaction', 'JobSatisfaction', 'MaritalStatus', 'MonthlyIncome', 'NumCompaniesWorked', 'WorkLifeBalance', 'YearsAtCompany'], dtype='object') y=dataframe['Attrition'] y.head() Out[36]: 0 0 1 2 1 Name: Attrition, dtype: int64 y, $x = dmatrices('Attrition \sim Age + Department + \$ DistanceFromHome + Education + EducationField + YearsAtCompany', dataframe, return_type="dataframe") print (x.columns) dtype='object') y = np.ravel(y)Train and Test data from sklearn.linear model import LogisticRegression model = LogisticRegression() model = model.fit(x, y)#Check the accuracy on the training set model.score(x, y) Out[39]: 0.8408163265306122 In [40]: y.mean() Out[40]: 0.16122448979591836 In [41]: X train, X test, y train, y test=sklearn.model selection.train test split(x,y, test size= model2=LogisticRegression() model2.fit(X train, y train) Out[41]: LogisticRegression() In [42]: model2.score(X train,y train) Out[42]: 0.8415937803692906 Training Accuracy of 84.15% is achieved by the model. In [43]: model2.score(X test,y test) Out[43]: 0.8435374149659864 Validation Accuracy of 84.35% is achieved by the model **Build and Evaluate Model** In [44]: predicted= model2.predict(X test) print (predicted) 0. 0. 0. 0. 0. 0. 0. 0.] In [45]: probs = model2.predict_proba(X_test) print(probs) [[0.86179622 0.13820378] [0.80754592 0.19245408] [0.74123931 0.25876069] [0.83441338 0.16558662] [0.73499935 0.26500065] [0.79097741 0.20902259] [0.85615196 0.14384804] [0.85699669 0.14300331] [0.96699058 0.03300942] [0.9368521 0.0631479] [0.95099278 0.04900722] [0.83101548 0.16898452] [0.86296556 0.13703444] [0.86581193 0.13418807] [0.88750604 0.11249396] [0.88892616 0.11107384] [0.88569727 0.11430273] [0.78516583 0.21483417] [0.7979449 0.2020551] [0.88511304 0.11488696] [0.70651589 0.29348411] [0.94676693 0.05323307] [0.86736254 0.13263746] [0.84276452 0.15723548] [0.60336837 0.39663163] [0.811292 0.188708 [0.91813731 0.08186269] [0.93285522 0.06714478] [0.68230751 0.31769249] [0.87027139 0.12972861] [0.87266386 0.12733614] [0.76968737 0.23031263] [0.86435774 0.13564226] [0.95758881 0.04241119] [0.84461487 0.15538513] [0.86719349 0.13280651] [0.90465982 0.09534018] [0.68936423 0.31063577] [0.90703618 0.09296382] [0.80663474 0.19336526] [0.91515727 0.08484273] [0.82351271 0.17648729] [0.93711517 0.06288483] [0.93411325 0.06588675] [0.89447655 0.10552345] [0.85317747 0.14682253] [0.78922387 0.21077613] [0.84879888 0.15120112] [0.66402447 0.33597553] [0.76252291 0.23747709] [0.92851112 0.07148888] [0.78953697 0.21046303] [0.86166592 0.13833408] [0.85837884 0.14162116] [0.87217674 0.12782326] [0.78950896 0.21049104] [0.87690794 0.12309206] [0.8416545 0.1583455] [0.72847141 0.27152859] [0.83181403 0.16818597] [0.90095034 0.09904966] [0.71077325 0.28922675] [0.92823024 0.07176976] [0.84375679 0.15624321] [0.79544106 0.20455894] [0.86826163 0.13173837] [0.91679452 0.08320548] [0.84763056 0.15236944] [0.89253708 0.10746292] [0.62872105 0.37127895] [0.93875394 0.06124606] [0.72620329 0.27379671] [0.85652974 0.14347026] [0.84226023 0.15773977] [0.77436384 0.22563616] [0.71899555 0.28100445] [0.9358739 0.0641261] [0.95710071 0.04289929] [0.79185834 0.20814166] [0.89370439 0.10629561] [0.9138204 0.0861796] [0.7935459 0.2064541] [0.77934019 0.22065981] [0.79638982 0.20361018] [0.83800499 0.16199501] [0.71395664 0.28604336] [0.97772717 0.02227283] [0.94645975 0.05354025] [0.88617626 0.11382374] [0.79620165 0.20379835] [0.61863823 0.38136177] [0.8186647 0.1813353] [0.74504124 0.25495876] [0.86779495 0.13220505] [0.8707114 0.1292886] [0.81717467 0.18282533] [0.71840765 0.28159235] [0.59825882 0.40174118] [0.83951552 0.16048448] [0.88351326 0.11648674] [0.74352571 0.25647429] [0.76631618 0.23368382] [0.98033038 0.01966962] [0.91857469 0.08142531] [0.7743284 0.2256716] [0.92514817 0.07485183] [0.88123384 0.11876616] [0.7458717 0.2541283] [0.90478362 0.09521638] [0.78685517 0.21314483] [0.81147771 0.18852229] [0.93472172 0.06527828] [0.93836504 0.06163496] [0.79411744 0.20588256] [0.813729 0.186271 [0.9161092 0.0838908] [0.90428346 0.09571654] [0.84669421 0.15330579] [0.95384554 0.04615446] [0.91283692 0.08716308] [0.85919605 0.14080395] [0.85902497 0.14097503] [0.87519517 0.12480483] [0.76114666 0.23885334] [0.92217687 0.07782313] [0.96859411 0.03140589] [0.94398222 0.05601778] [0.81780289 0.18219711] [0.88058705 0.11941295] [0.77894286 0.22105714] [0.97124466 0.02875534] [0.88807662 0.11192338] [0.78715262 0.21284738] [0.8200148 0.1799852] [0.94934545 0.05065455] [0.95888935 0.04111065] [0.73559214 0.26440786] [0.93416998 0.06583002] [0.73750627 0.26249373] [0.82136746 0.17863254] [0.821712 0.178288] [0.89896704 0.10103296] [0.78745764 0.21254236] [0.8982535 0.1017465] [0.9143382 0.0856618] [0.92724744 0.07275256] [0.96594967 0.03405033] [0.94417368 0.05582632] [0.93073083 0.06926917] [0.66320576 0.33679424]

[0.84168649 0.15831351]

