



Team2: HongSeop, Lynda, Roger, Yashasv, Megumu



DATA AND OBJECTIVE



- IBM HR Attrition Dataset
- The dataset has 35 variable columns and 1,470 rows

ge	- Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	Edu	ıcationField	EmployeeCount	EmployeeNumber	EnvironmentSatisfaction		Gender	HourlyRate	Jobinvolvement	JobLevel	
igint	string	string	bigint	string	tinyint	tinyint	stri	ing	tinyint	bigint	tinyint		string	int	tinyint	tinyint	
nteger	Boolean	Text	Integer	Text	Integer	Integer	Text		Integer	Integer	Integer		Gender	Integer	Integer	Integer	
4	1 Yes	Travel_Rarely	1102	Sales	1		2 Life	Sciences	1	1 1		2	Female	94		3	2
4	9 No	Travel_Frequently	279	Research & Development	8		1 Life	Sciences	1	1 2		3	Male	61		2	2
3	7 Yes	Travel_Rarely	1373	Research & Development	2		2 Othe	er	1	1 4		4	Male	92	2	2	1
3	3 No	Travel_Frequently	1392	Research & Development	3		4 Life	Sciences	1	1 5		4	Female	56	5	3	1
2	7 No	Travel_Rarely	591	Research & Development	2		1 Med	dical	1	1 7		1	Male	40		3	1
3	2 No	Travel_Frequently	1005	Research & Development	2		2 Life	Sciences	1	1 8		4	Male	79		3	1
5	9 No	Travel_Rarely	1324	Research & Development	3		3 Med	dical	1	1 10		3	Female	81		4	1
3	0 No	Travel_Rarely	1358	Research & Development	24		1 Life:	Sciences	1	1 11		4	Male	67		3	1
3	8 No	Travel_Frequently	216	Research & Development	23		3 Life	Sciences	1	1 12		4	Male	44		2	3
3	6 No	Travel_Rarely	1299	Research & Development	27		3 Med	dical	1	1 13		3	Male	94		3	2
3	5 No	Travel_Rarely	809	Research & Development	16		3 Med	dical	1	1 14		1	Male	84	1	4	1
2	9 No	Travel_Rarely	153	Research & Development	15		2 Life	Sciences	1	1 15		4	Female	49)	2	2
3	1 No	Travel_Rarely	670	Research & Development	26		1 Life:	Sciences	1	1 16		1	Male	31		3	1

Objective

How to analyze

How the company decrease leaving employees

Finding features(columns) related to Attrition: "Yes"



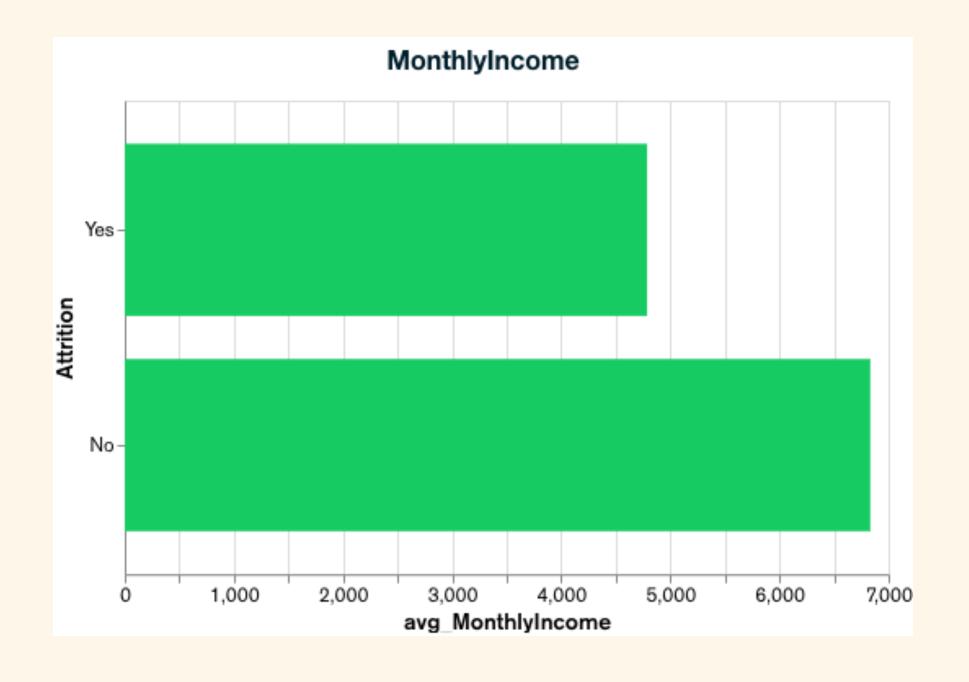
DATA ANALYZING FLOW

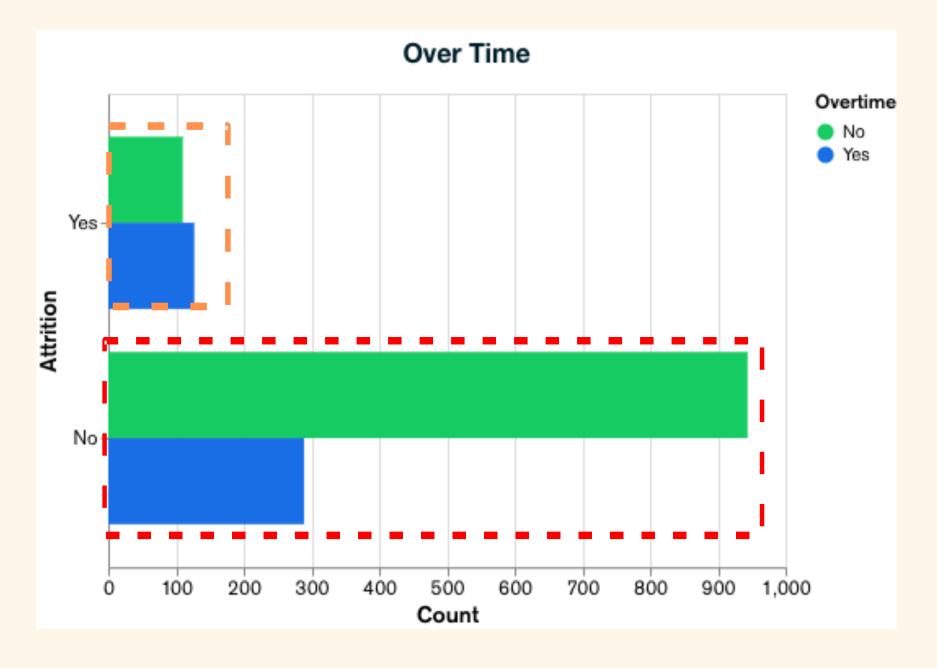
	Data Visualizing	Al Analysis	Business Insight Recommendation
Contents	 Visualizing the data Finding the factor related to Attrition on the visualized charts 	 Making the AI model Checking the coefficients to confirm the factors related to Attrition or not 	 Gaining the business insights based on our visualizing and AI analysis Recommendations to IBM
DB	Mongo DB	Spark	-
Environment	Mongo DB Atlas Dashboard	PySpark on Dataiku	-



DATA VISUALIZING IN MONGODB

We visualized almost of all data. In these data, average of "MonthlyIncome" and "OverTime" has a big difference between Attrition YES and No.







PREDICTIVE MODELING IN DFS -**PYSPARK**

Al prediction results

+	+		+	+		
features	label	rawPrediction	probability	prediction		
+	+		+	+		
(5,[0,3],[1.0,246	0	[1.95241127441409	[0.87570932879953	0.0		
(5,[0,3],[1.0,254	0	[1.96248345523969	[0.87680146802390	0.0		
(5,[0,3],[1.0,261	1	[1.97123034911454	[0.87774320340243	0.0		
(5,[0,3],[1.0,261	0	[1.97162793519976	[0.87778586198671	0.0		
(5,[0,3],[1.0,268	0	[1.98103747254999	[0.87879171383867	0.0		
(5,[0,3],[1.0,294	0	[2.01522987587897	[0.88238686655610	0.0		
(5,[0,3],[1.0,455	0	[2.22913118972775	[0.90283517032261	0.0		
(5,[0,3],[1.0,472	0	[2.25139601050011	[0.90477088409279	0.0		
(5,[0,3],[1.0,506	0	[2.29632323813006	[0.90857207303412	0.0		
(5,[0,3],[1.0,513	0	[2.30533518939506	[0.90931793422989	0.0		
(5,[0,3],[1.0,525	0	[2.32203380497434	[0.91068550334029	0.0		
(5,[0,3],[1.0,560	0	[2.36815379085995	[0.91436641205485	0.0		
(5,[0,3],[1.0,596	0	[2.41626170717166	[0.91805896423421	0.0		
(5,[0,3],[1.0,639	0	[2.47311651735823	[0.92223556685983	0.0		
(5,[0,3],[1.0,894	0	[2.81053457501559	[0.94324244505347	0.0		
(5,[0,3],[1.0,121	1	[3.23807214532299	[0.96224212864583	0.0		
(5,[0,3],[1.0,138	0	[3.45767219305993	[0.96945911998152	0.0		
(5,[0,3],[1.0,152	0	[3.6400316774812,	[0.97442000128206	0.0		
(5,[0,3],[1.0,178	0	[3.99176283420650	[0.98186772026493	0.0		
(5,[0,3],[1.0,184	0	[4.06783430517875	[0.98317356161897	0.0		
+	+		+	+		
only showing top 20 rows						

AUC ROC:0.8008662508662475

Coefficients

No	Features	Feature Weight	Note		
1	BTFreq	0.473	"BusinessTravel" Travel-Frequently: 1 Travel-Rarely: 0 No-Travel: 0		
2	MSSingle	0.899	"MaritalStatus" Single:1 Discovered:0 Married:0		
3	OverTime	1.190	row data		
4	MonthlyIncome	-1.115	row data		
5	Department	0.561	"Department" Sales:1 HR:1 Researcher:0		

Model intercept: -2.124



PREDICTIVE MODELING INSIGHTS



AUC ROC = 0.80 -> Predicts Y variables with high confidence



Single, Overtime -> Likely to leave



Monthly Income -> Likely to stay



Overall, it aligns with the previous MongoDB analysis



CONCLUSIONS

- Overtime work has a significant direct relationship with Attrition: the more overtime employees work, the more likely they leave IBM.
- Salary increases and bonuses are not features that significantly impact Attrition, although there is always room for improvement, this can slightly reduce the chances of leaving.
- Mini vacations in between trips and organizing in-house social networking events reduce the odds of Attrition=yes.
- Offer career development opportunities has a relationship with Attrition and reduces the odds that employees -especially singles- leave IBM.



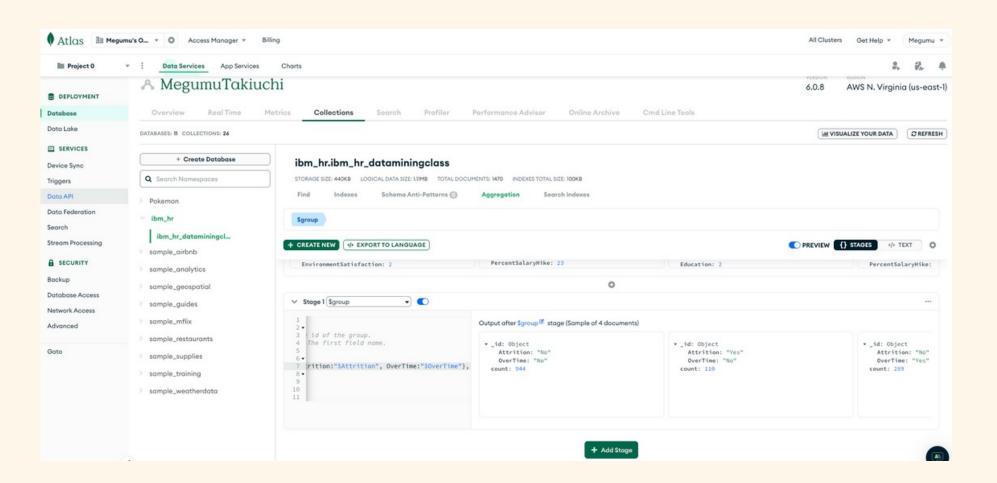
REFERENCES

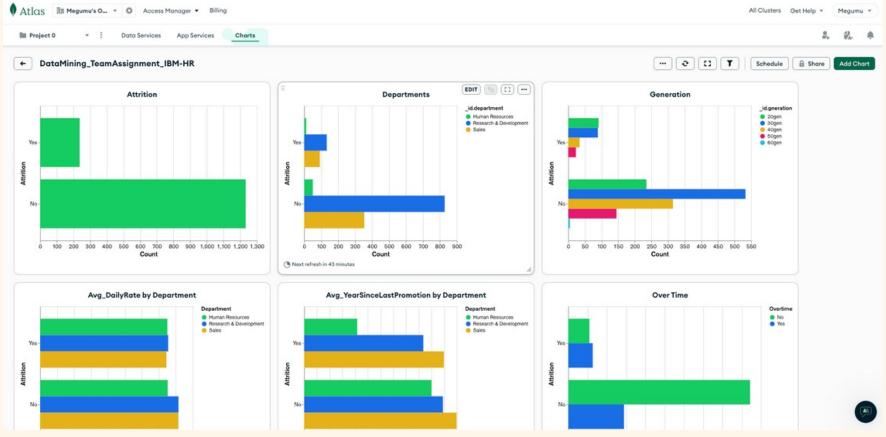
Bock, L. (2022, July). It's Time to Reimagine Employee Retention. Harvard Business Review. Retrieved from https://hbr.org/2022/07/its-time-to-reimagine-employee-retention



APPENDIX

- We used MongoDB and Spark to practice visualizing and creating AI models.
- Visualizing environment: Mongo DB Atlas







PySpark Code (1)

```
from pyspark.ml import Pipeline
     from pyspark.ml.classification import LogisticRegression
    from pyspark.ml.feature import HashingTF, Tokenizer
     from pyspark.sql import Row
     from pyspark.sql.functions import UserDefinedFunction
     from pyspark.sql.types import *
     from pyspark.ml.feature import VectorAssembler
     from pyspark.mllib.evaluation import *
     #creating vectors with names of variables
10
     vecAssembler = VectorAssembler(inputCols = ['BTFreq', 'MSSingle', 'OverTime', 'MonthlyIncome', 'HRSales'], output
11
    v df = vecAssembler.transform(df)
12
     vhouse df = v df.select(['features', 'Attrition'])
13
     vhouse df = vhouse df.withColumnRenamed("Attrition", "label")# We have to rename our output variable to 'label'
15
16
     #splitting the dataset
     splits = vhouse df.randomSplit([0.7, 0.3])
17
    train df = splits[0]
     test df = splits[1]
20
     #creating an object with the logistic regression engine
21
     lr = LogisticRegression(maxIter=20)
     pipeline = Pipeline(stages=[lr])
23
24
25
     #fitting the model
     model = lr.fit(train df)
26
27
     #evaluating the model using testing data
     result = model.transform(test_df)
     result.prediction
30
     result.show()
31
32
     from pyspark.ml.evaluation import BinaryClassificationEvaluator
33
     evaluator = BinaryClassificationEvaluator(rawPredictionCol="rawPrediction")
34
     AUC ROC = evaluator.evaluate(result, {evaluator.metricName: "areaUnderROC"})
35
36
     print('AUC ROC:' + str(AUC ROC))
```



PySpark Code (2)

```
from pyspark.ml.evaluation import BinaryClassificationEvaluator
41
     # Evaluate model
42
     evaluator = BinaryClassificationEvaluator(rawPredictionCol="rawPrediction")
     evaluator.evaluate(result)
45
     from pyspark.ml.tuning import ParamGridBuilder, CrossValidator
47
     # Create ParamGrid for Cross Validation
49 ▼ paramGrid = (ParamGridBuilder()
50
                  .addGrid(lr.regParam, [0.01, 0.5, 2.0])
                  .addGrid(lr.elasticNetParam, [0.0, 0.5, 1.0])
51
52
                  .addGrid(lr.maxIter, [1, 5, 10])
53
                  .build())
    # Create 5-fold CrossValidator
     cv = CrossValidator(estimator=lr, estimatorParamMaps=paramGrid, evaluator=evaluator, numFolds=5)
56
57 # Run cross validations
    cvModel = cv.fit(train df)
    # this will likely take a fair amount of time because of the amount of models that we're creating and testing
60
    # Use test set to measure the accuracy of our model on new data
     predictions = cvModel.transform(test df)
63
    # cvModel uses the best model found from the Cross Validation
     # Evaluate best model
     evaluator.evaluate(predictions)
67
    print('Model Intercept: ', cvModel.bestModel.intercept)
     weights = cvModel.bestModel.coefficients
     weights = [(float(w),) for w in weights] # convert numpy type to float, and to tuple
    weightsDF = sqlContext.createDataFrame(weights, ["Feature Weight"])
    weightsDF.show()
73 # View best model's predictions and probabilities of each prediction class
     selected = predictions.select("label", "prediction", "probability", "features")
75 selected.show()
```



