

SC1015 Mini Project

DSF3
HUNG KUO-CHEN
JODI TAY SEOW XUAN
YANG XIAOYUE

Motivation

- Employee attrition: the turnover rate in various job roles
- High attrition rate brings about problems
 - Hard-to-replace employees leave → lower productivity and profits
 - High costs incurred in training and hiring new employees
- Aim: uncover reasons for an employee's resignation and recommend improvements made within IBM to retain its employees



IBM HR Attrition Dataset

- Used data from 1470 employees in IBM
 - 16.1% left IBM, 83.9% stayed in IBM
- Includes 34 independent variables based on an employee's profile
 - Each contributed to whether an employee decided to leave or stay



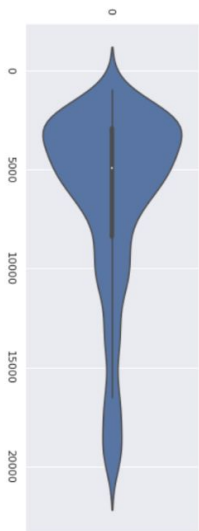
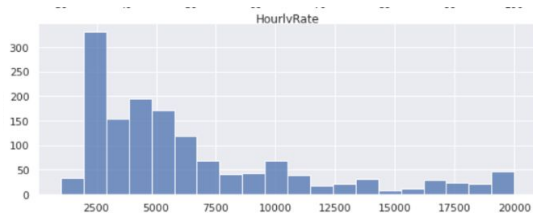
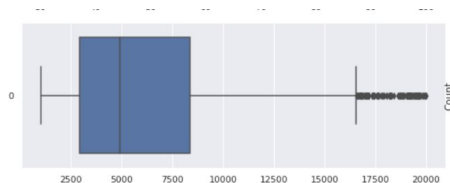
Cleaning the Data

- Dropped insignificant columns
 - EmployeeCount, Over18, StandardHours, EmployeeNumber
- Checked for missing values and duplicates
 - None found

Breakdown of Variables

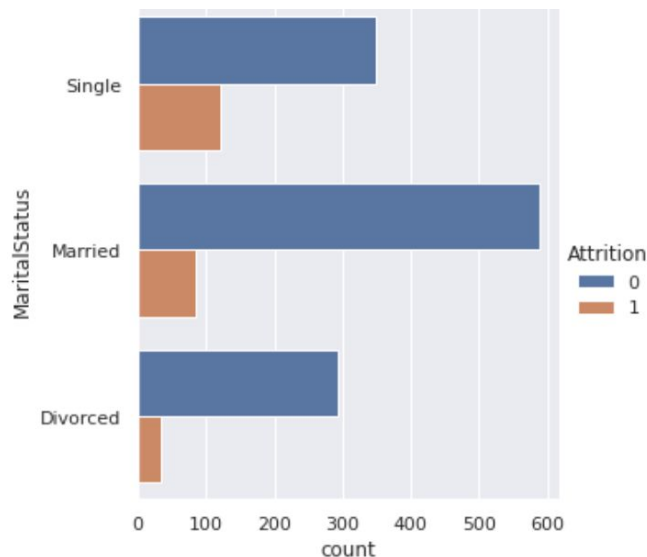
Numeric

- 14 regular numeric variables
- 9 factor numeric variables
- A variety of data visualization methods: box plot, histogram plot, and violin plot



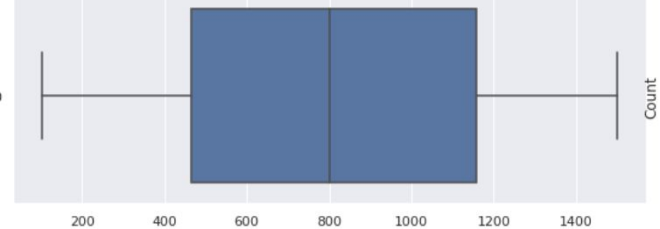
Categorical

- 6 categorical variables
- Categorical bar plot of each variable against attrition using GroupBy

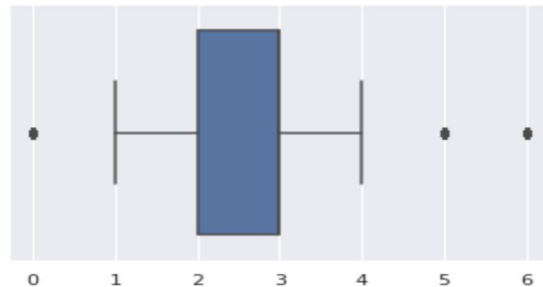


Outliers (Box Plot)

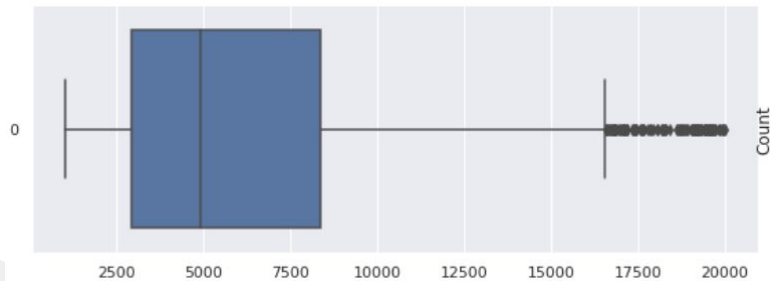
Different variables have different degrees of outliers (close to none, moderate, and large)



Box plot with close to no outliers



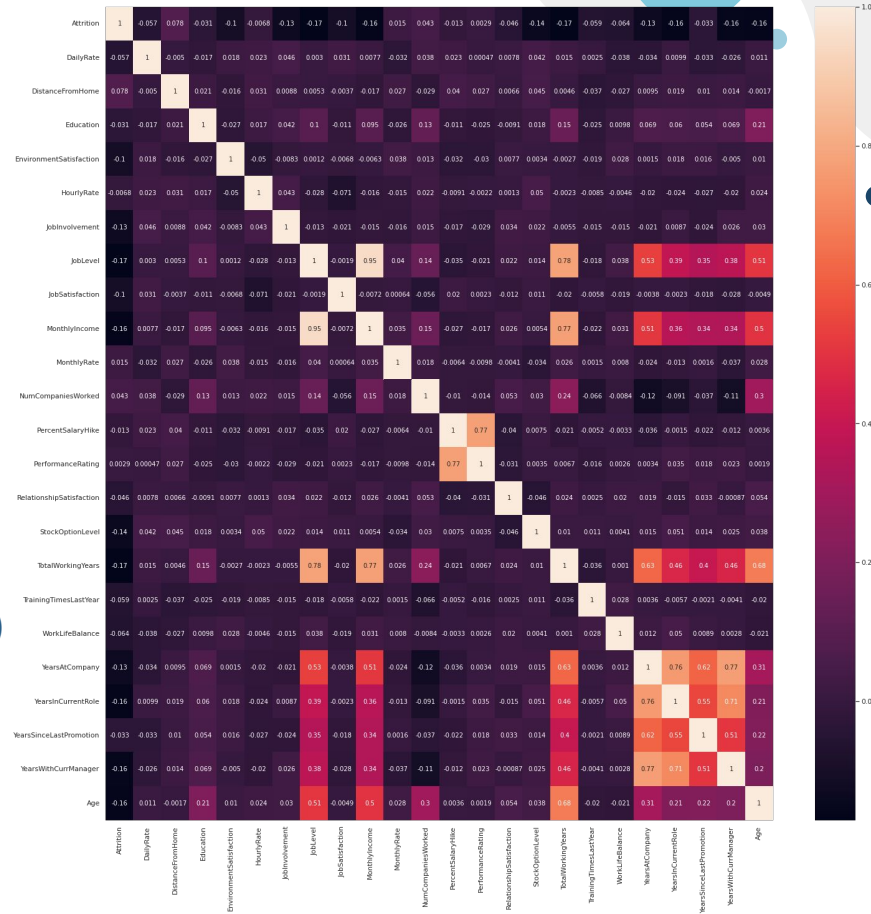
Box plot with a moderate number of outliers



Box plot with a large number of outliers

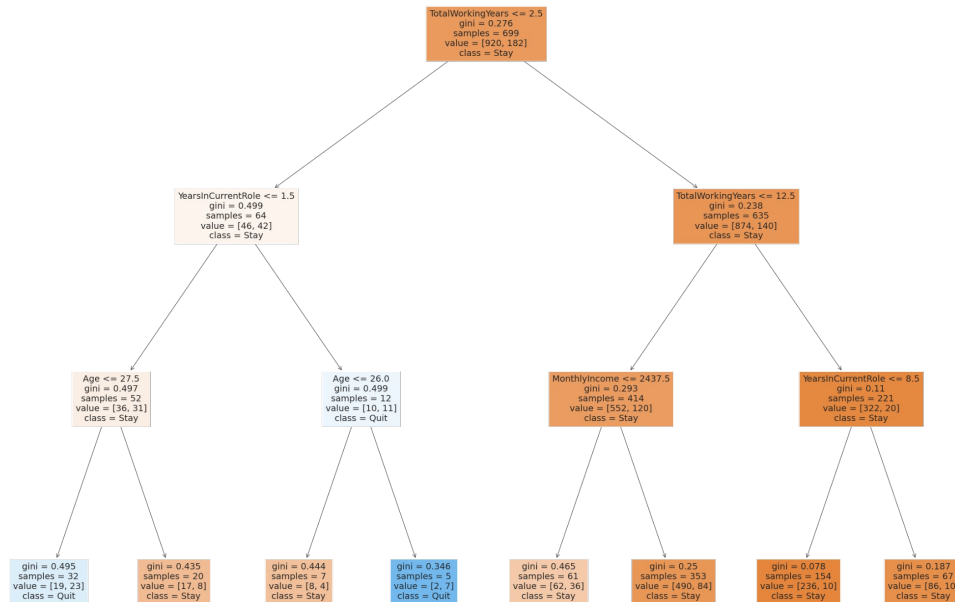
Correlation Matrix

- Reclassify attrition as a numeric variable
- Plot the correlation matrix for attrition against all other numeric variables
- Interesting Findings:
 - MonthlyIncome and JobLevel (0.95)
 - TotalWorkingYears and JobLevel (0.78)
 - TotalWorkingYears and MonthlyIncome (0.77)



Random Forest

- Extract variables with relatively high correlation
- **Accuracy of random forest: 84.5%**
- Tune hyperparameters with random search
 - Random combinations of the hyperparameters are used to find the optimal solution for the built model



Logistic Regression

- Determine level of influence of categorical variables on attrition
- Convert each variable into numeric indicator variables with `get_dummies`
- Accuracy: 0.84 (train), 0.86 (test)
- Ineffective due to **extremely low precision**
 - The model only classifies 38% of employees that quit correctly

	OverTime	Gender	MaritalStatus	Department	EducationField
0	Yes	Female	Single	Sales	Life Sciences
1	No	Male	Married	Research & Development	Life Sciences
2	Yes	Male	Single	Research & Development	Other
3	Yes	Female	Married	Research & Development	Life Sciences
4	No	Male	Married	Research & Development	Medical
...
1465	No	Male	Married	Research & Development	Medical
1466	No	Male	Married	Research & Development	Medical
1467	Yes	Male	Married	Research & Development	Life Sciences
1468	No	Male	Married	Sales	Medical
1469	No	Male	Married	Research & Development	Medical

1470 rows x 5 columns

Neural Network

- Train a simple multilayer perceptron model by three numeric attributes which is highly related to attrition
- The **loss** of the model **reduced strikingly** after three epochs
- Accuracy: 83.4% (train), **85.7% (test)**

```
Epoch 1/3, Iteration 1/12, Loss: 641.4485
Epoch 1/3, Iteration 2/12, Loss: 2409.5425
Epoch 1/3, Iteration 3/12, Loss: 641.4485
Epoch 1/3, Iteration 4/12, Loss: 2409.5425
Epoch 1/3, Iteration 5/12, Loss: 641.4485
Epoch 1/3, Iteration 6/12, Loss: 2409.5425
Epoch 1/3, Iteration 7/12, Loss: 641.4485
Epoch 1/3, Iteration 8/12, Loss: 2409.5425
Epoch 1/3, Iteration 9/12, Loss: 641.4485
Epoch 1/3, Iteration 10/12, Loss: 2409.5425
Epoch 1/3, Iteration 11/12, Loss: 641.4485
Epoch 1/3, Iteration 12/12, Loss: 2409.5425
Epoch 2/3, Iteration 1/12, Loss: 0.4801
Epoch 2/3, Iteration 2/12, Loss: 0.4577
Epoch 2/3, Iteration 3/12, Loss: 0.3911
Epoch 2/3, Iteration 4/12, Loss: 0.4347
Epoch 2/3, Iteration 5/12, Loss: 0.4776
Epoch 2/3, Iteration 6/12, Loss: 0.4347
Epoch 2/3, Iteration 7/12, Loss: 0.4776
Epoch 2/3, Iteration 8/12, Loss: 0.4347
Epoch 2/3, Iteration 9/12, Loss: 0.4776
Epoch 2/3, Iteration 10/12, Loss: 0.4347
Epoch 2/3, Iteration 11/12, Loss: 0.4776
Epoch 2/3, Iteration 12/12, Loss: 0.4347
Epoch 3/3, Iteration 1/12, Loss: 0.4801
Epoch 3/3, Iteration 2/12, Loss: 0.4577
Epoch 3/3, Iteration 3/12, Loss: 0.3911
Epoch 3/3, Iteration 4/12, Loss: 0.4347
Epoch 3/3, Iteration 5/12, Loss: 0.4776
Epoch 3/3, Iteration 6/12, Loss: 0.4347
Epoch 3/3, Iteration 7/12, Loss: 0.4776
Epoch 3/3, Iteration 8/12, Loss: 0.4347
Epoch 3/3, Iteration 9/12, Loss: 0.4776
Epoch 3/3, Iteration 10/12, Loss: 0.4347
Epoch 3/3, Iteration 11/12, Loss: 0.4776
Epoch 3/3, Iteration 12/12, Loss: 0.4347
```

```
Net(
  (fc1): Linear(in_features=3, out_features=100, bias=True)
  (relu1): ReLU()
  (fc2): Linear(in_features=100, out_features=50, bias=True)
  (relu2): ReLU()
  (fc3): Linear(in_features=50, out_features=2, bias=True)
)
```

Epoch	Iteration	Loss	Accuracy
1	1	641.4485	0.0000
1	2	2409.5425	0.0000
1	3	641.4485	0.0000
1	4	2409.5425	0.0000
1	5	641.4485	0.0000
1	6	2409.5425	0.0000
1	7	641.4485	0.0000
1	8	2409.5425	0.0000
1	9	641.4485	0.0000
1	10	2409.5425	0.0000
1	11	641.4485	0.0000
1	12	2409.5425	0.0000
2	1	0.4801	0.0000
2	2	0.4577	0.0000
2	3	0.3911	0.0000
2	4	0.4347	0.0000
2	5	0.4776	0.0000
2	6	0.4347	0.0000
2	7	0.4776	0.0000
2	8	0.4347	0.0000
2	9	0.4776	0.0000
2	10	0.4347	0.0000
2	11	0.4776	0.0000
2	12	0.4347	0.0000
3	1	0.4801	0.0000
3	2	0.4577	0.0000
3	3	0.3911	0.0000
3	4	0.4347	0.0000
3	5	0.4776	0.0000
3	6	0.4347	0.0000
3	7	0.4776	0.0000
3	8	0.4347	0.0000
3	9	0.4776	0.0000
3	10	0.4347	0.0000
3	11	0.4776	0.0000
3	12	0.4347	0.0000

```
# Out loss function
criterion = nn.CrossEntropyLoss()

# Our optimizer
learning_rate = 0.0001
optimizer = torch.optim.SGD(net.parameters(), lr=learning_rate, nesterov=True, momentum=0.9, dampening=0)
```



Summary of Findings

- 1. **Random Forest**

Numeric variables that have relatively high correlation with attrition can be used to predict attrition

- 2. **Logistic Regression**

Not recommended using only categorical variables to predict attrition

- 3. **Neural Network**

The trained model can be implemented to predict attrition effectively

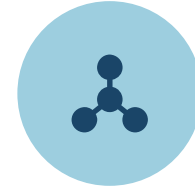


Data Driven Insights



Reasons for Leaving

- Low salary
- Low chance for career progression
- Lack of opportunities
- Long distance from residence



Improvements for IBM

- Provide more salary incentives or other allowance
- Enhance effective employee assessments
- Open up spots for changes in senior management

Thank You

CREDITS: This presentation template was created by **Slidesgo**, including icons by **Flaticon**, infographics & images by **Freepik**