



# Employee Churn Analysis

USING MACHINE LEARNING

DMPA LAB PROJECT

- Analysis of dataset

## Overview

Overview	Reproduction	Warnings 1
Dataset statistics		
Number of variables	10	
Number of observations	14999	
Missing cells	0	
Missing cells (%)	0.0%	
Duplicate rows	3008	
Duplicate rows (%)	20.1%	
Total size in memory	1.1 MIB	
Average record size in memory	80.0 B	
Variable types		
NUM	5	
BOOL	3	
CAT	2	

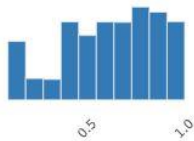
## Variables

satisfaction\_level

Real number ( $\mathbb{R}_{\geq 0}$ )

Distinct count	92
Unique (%)	0.6%
Missing	0
Missing (%)	0.0%
Infinite	0
Infinite (%)	0.0%

Mean	0.6128335222348156
Minimum	0.09
Maximum	1.0
Zeros	0
Zeros (%)	0.0%
Memory size	117.2 KiB



A histogram showing the distribution of satisfaction\_level values. The x-axis ranges from 0.5 to 1.0, and the y-axis represents frequency. The distribution is unimodal and slightly right-skewed, with a peak around 0.75.

Bin Range	Frequency
0.45 - 0.50	10
0.50 - 0.55	5
0.55 - 0.60	5
0.60 - 0.65	15
0.65 - 0.70	12
0.70 - 0.75	13
0.75 - 0.80	14
0.80 - 0.85	15
0.85 - 0.90	16
0.90 - 0.95	14

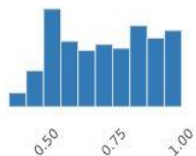
Toggle details

last\_evaluation

Real number ( $\mathbb{R}_{\geq 0}$ )

Distinct count	65
Unique (%)	0.4%
Missing	0
Missing (%)	0.0%
Infinite	0
Infinite (%)	0.0%

Mean	0.7161017401160078
Minimum	0.36
Maximum	1.0
Zeros	0
Zeros (%)	0.0%
Memory size	117.2 KiB



A histogram showing the distribution of 'last\_evaluation' values. The x-axis is labeled with 0.50, 0.75, and 1.00. The y-axis represents frequency. The distribution is roughly bell-shaped, centered around 0.75, with a peak frequency of approximately 10 for the bin [0.65, 0.75].

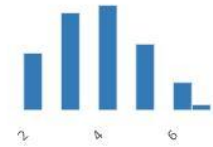
Toggle details

### number\_project

Real number ( $\mathbb{R}_{\geq 0}$ )

Distinct count	6
Unique (%)	< 0.1%
Missing	0
Missing (%)	0.0%
Infinite	0
Infinite (%)	0.0%

Mean	3.80305353690246
Minimum	2
Maximum	7
Zeros	0
Zeros (%)	0.0%
Memory size	117.2 KiB



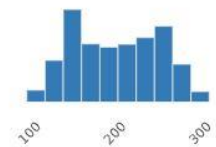
Toggle details

### average\_monthly\_hours

Real number ( $\mathbb{R}_{\geq 0}$ )

Distinct count	215
Unique (%)	1.4%
Missing	0
Missing (%)	0.0%
Infinite	0
Infinite (%)	0.0%

Mean	201.0503366891126
Minimum	96
Maximum	310
Zeros	0
Zeros (%)	0.0%
Memory size	117.2 KiB



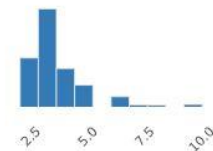
Toggle details

### time\_spend\_company

Real number ( $\mathbb{R}_{\geq 0}$ )

Distinct count	8
Unique (%)	0.1%
Missing	0
Missing (%)	0.0%
Infinite	0
Infinite (%)	0.0%

Mean	3.498233215547703
Minimum	2
Maximum	10
Zeros	0
Zeros (%)	0.0%
Memory size	117.2 KiB



Toggle details

### Work\_accident

Boolean

Distinct count	2
Unique (%)	< 0.1%
Missing	0
Missing (%)	0.0%
Memory size	117.2 KiB



Toggle details

### promotion\_last\_5years

Boolean

Distinct count	2
Unique (%)	< 0.1%
Missing	0
Missing (%)	0.0%
Memory size	117.2 KiB

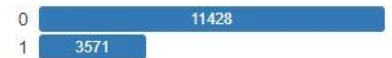


[Toggle details](#)

### quit

Boolean

Distinct count	2
Unique (%)	< 0.1%
Missing	0
Missing (%)	0.0%
Memory size	117.2 KiB

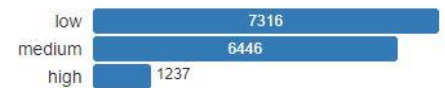


[Toggle details](#)

### salary

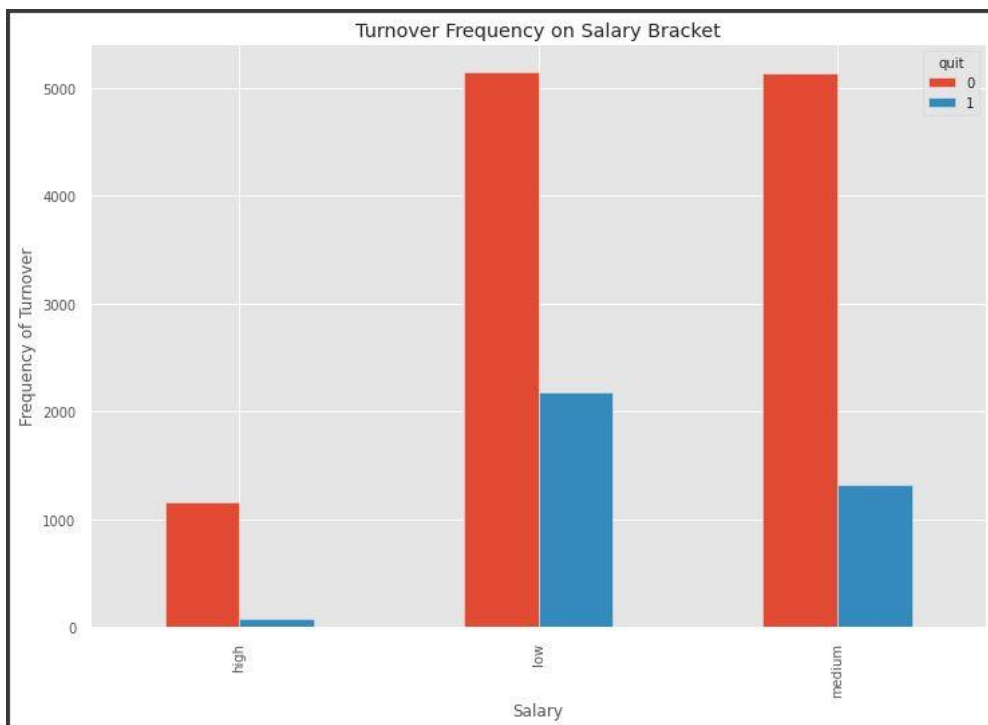
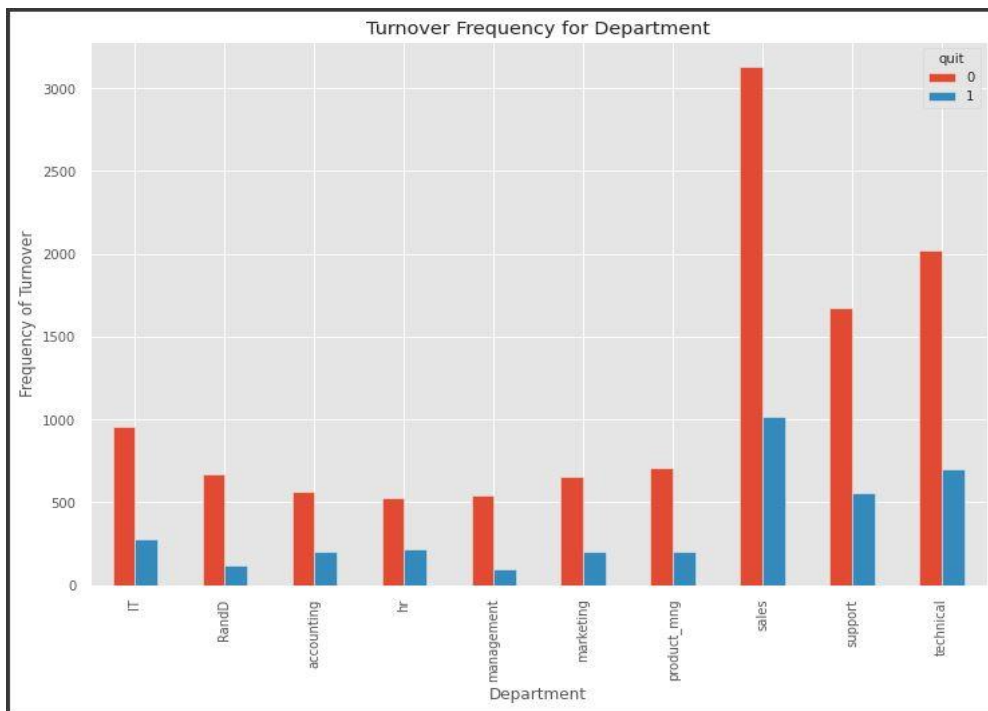
Categorical

Distinct count	3
Unique (%)	< 0.1%
Missing	0
Missing (%)	0.0%
Memory size	117.2 KiB



[Toggle details](#)

- ANALYSIS OF CATEGORICAL FEATURES



- COMPARISON OF ALGORITHMS USED

1. Decision Tree

```
from sklearn.tree import DecisionTreeClassifier
from sklearn import metrics

# Training the Decision Tree Classification model on the Training set
model=DecisionTreeClassifier()
model=model.fit(X_train,y_train)

# Prediction
y_pred=model.predict(X_test)

# Accuracy
print("\nAccuracy (Decision Tree Classifier Model):",metrics.accuracy_score(y_test, y_pred))
```

Accuracy (Decision Tree Classifier Model): 0.981

crit

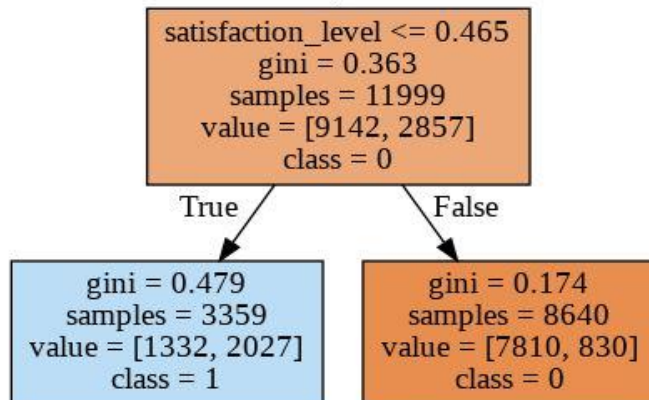
split

depth

min\_split

min\_leaf

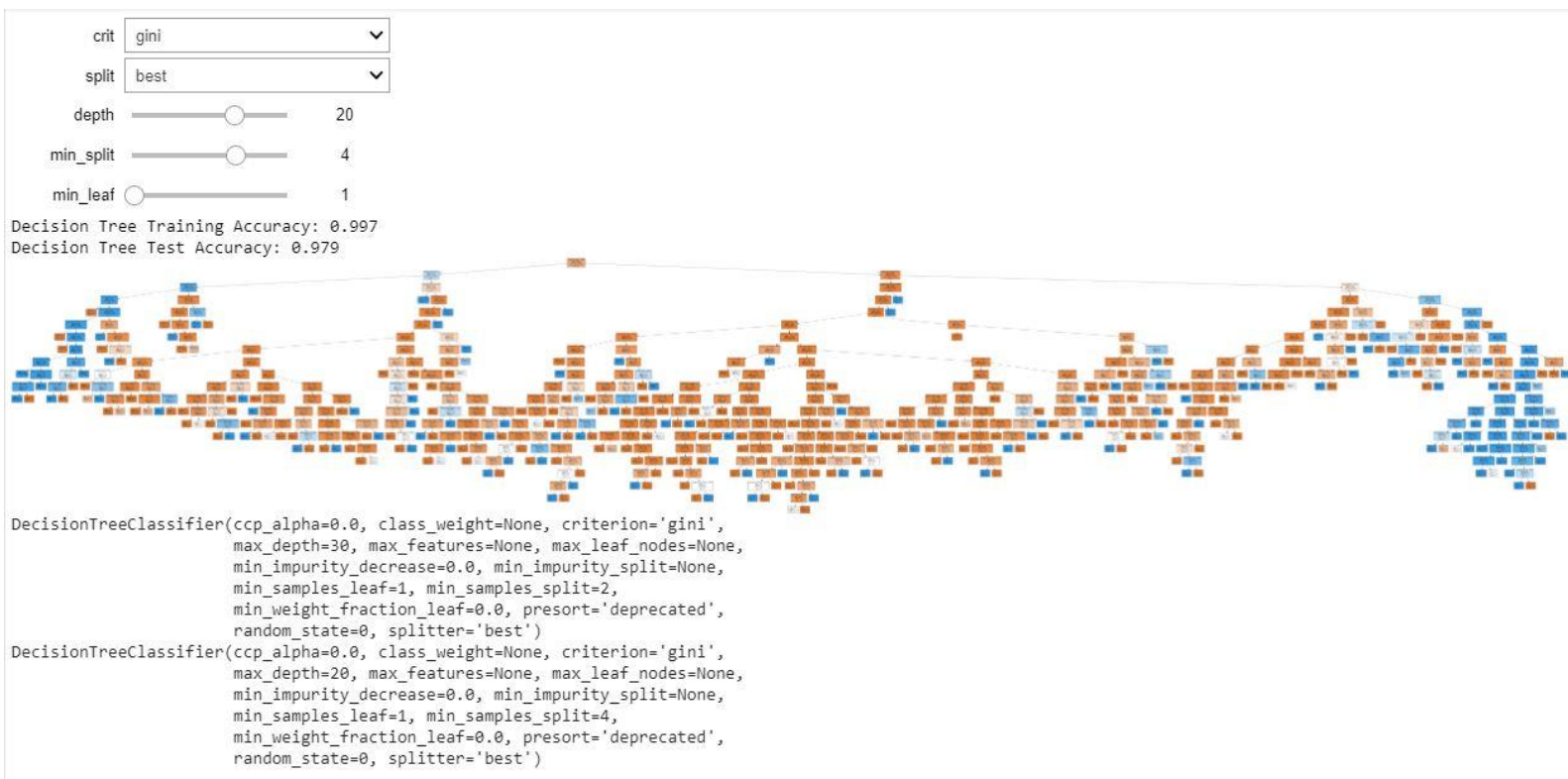
Decision Tree Training Accuracy: 0.820  
Decision Tree Test Accuracy: 0.823



```
DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='gini',
max_depth=1, max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, presort='deprecated',
random_state=0, splitter='best')
```

Widget for analysis of Decision Tree





## 2. Random Forest

[21] # Training the Random Forest Classification model on the Training set

```
from sklearn.ensemble import RandomForestClassifier
classifier = RandomForestClassifier(n_estimators = 10, criterion = 'entropy', random_state = 0)
classifier.fit(X_train, y_train)

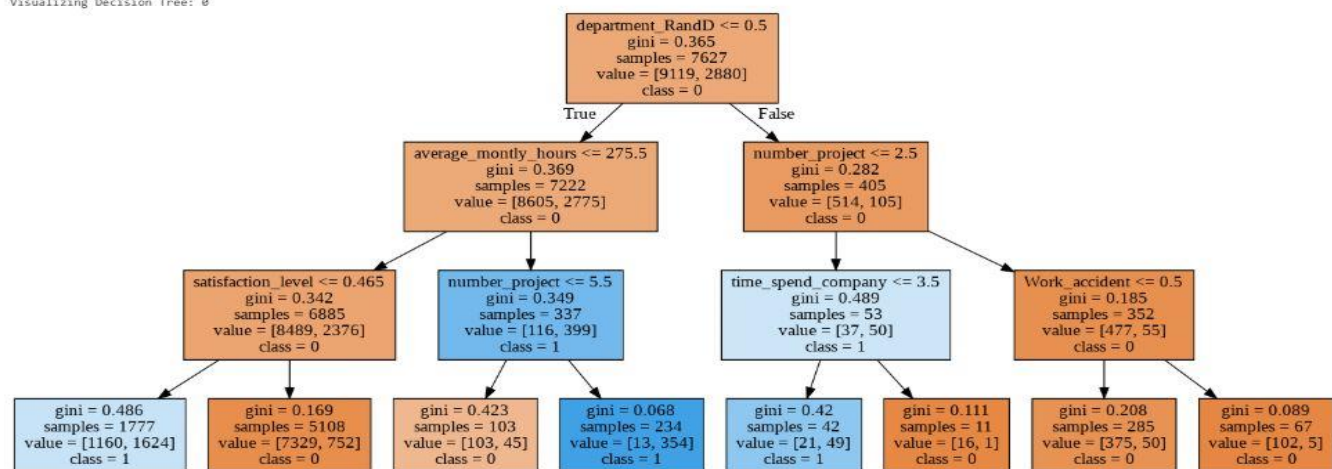
# Prediction
y_pred = classifier.predict(X_test)

# Accuracy
print("\nAccuracy (Random Forest Classifier Model):",metrics.accuracy_score(y_test, y_pred))
```



Accuracy (Random Forest Classifier Model): 0.988

[22] ☐    
 depth   
 forests   
 min\_split   
 min\_leaf   
 Random Forest Training Accuracy: 0.919  
 Random Forest Test Accuracy: 0.917  
 Visualizing Decision Tree: 0



RandomForestClassifier(bootstrap='True', ccp\_alpha=0.0, class\_weight=None, criterion='gini', max\_depth=3, max\_features='auto', max\_leaf\_nodes=None, max\_samples=None, min\_impurity\_decrease=0.0, min\_impurity\_split=None, min\_samples\_leaf=1, min\_samples\_split=2, min\_weight\_fraction\_leaf=0.0, n\_estimators=100, n\_iter\_no\_change=None, random\_state=None)

### Widget for analysis of Random Forest

### 3. KNN

```
[26] # Training the K-NN model on the Training set
from sklearn.neighbors import KNeighborsClassifier
classifier = KNeighborsClassifier(n_neighbors = 5, metric = 'minkowski', p = 2)
classifier.fit(X_train, y_train)

# Prediction
y_pred = classifier.predict(X_test)

# Accuracy
print("\nAccuracy (Random Forest Classifier Model):", metrics.accuracy_score(y_test, y_pred))
```



Accuracy (Random Forest Classifier Model): 0.9326666666666666



#### 4. Logistic Regression

```
[24] from sklearn.linear_model import LogisticRegressionCV

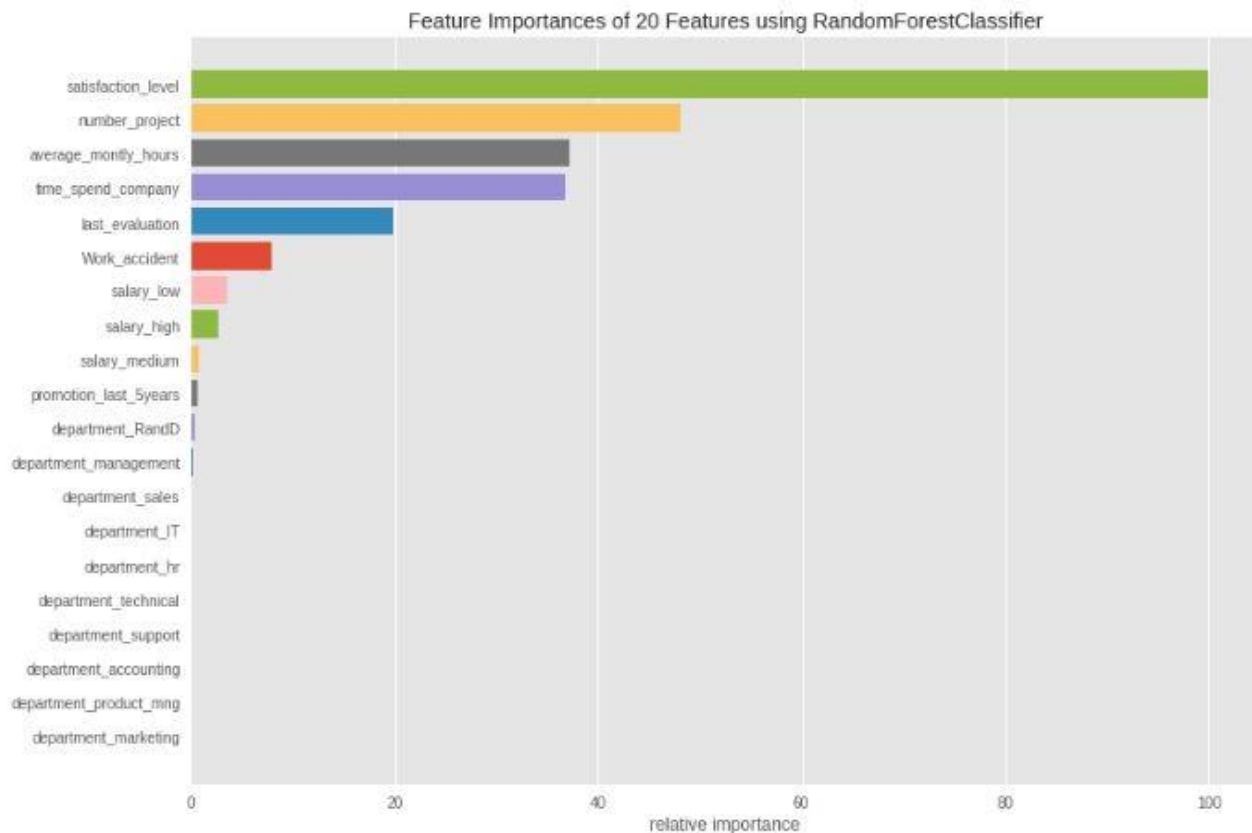
logit = LogisticRegressionCV(random_state=1, n_jobs=-1,max_iter=500,
                             cv=10)

lr = logit.fit(X_train, y_train)

print('Logistic Regression Accuracy: {:.3f}'.format(accuracy_score(y_test, lr.predict(X_test))))
```

Logistic Regression Accuracy: 0.790

- Relative Importance of Features



For complete analysis, graphs and code visit the following link:

<https://github.com/apoorv-mit2021/Employee-Churn-Analysis>