

# Predict Employee Turnover with SciKit Learn

- It is a free software machine learning library for Python programming language.
- Predict if employees are likely to quit or not.
- Quit → 1 means yes.

## TASK-1 IMPORT LIBRARY

## TASK-2 EXPLORATORY DATA ANALYSIS

[ `hr = pd.read_csv('data/employee_data.csv')` → Import  
`hr.head()` → first few entries

[ `hr.profile-report (title="Data")` → Overview report of data

[ `pd.crosstab(hr.salary, hr.quit).plot(kind='bar')`   
 `plt.title("Turnover freq on Salary Bracket")`   
 `plt.xlabel('Salary')`   
 `plt.ylabel('Frequency of Turnover')`   
 `plt.show()`   
 *axis axis bar graph*

## TASK-3 Encode Categorical Feature

→ The variables are categorical i.e. string or object but scikit-learn and many others can't process these so encoding is done.

We will use one ~~has~~<sup>hot</sup> encoding. Matlab  
jaise ek employee ke dept mein accounting  
hai toh table banega and use  
accounting ke age 1 baki ke age 0.

**DANGER**

pd.get\_dummies() → age direct use krenge toh  
data lose hoga.

category-variables = ['department', 'salary']  
for var in category-variables:  
    category-list = pd.get\_dummies  
        (hs[var], prefix=var)  
hs = hs.join(category-list)

name same  
keliye toh  
data ko lose  
ho

hs.head()

Ab hum main salary and department delete  
kridenge.

hs.drop(columns = ['department', 'salary'],  
        axis=1, inplace = True)

isko matlab vertical

PTO →



**TASK 4 - Visualize Class Imbalance** → Agar sab kisi ke hawafod testing mein fail ho rha ho toh ye problem ho rhi hai. In baad check kr lena chahiye.

```
from yellowbrick.target import ClassBalance
plt.style.use('ggplot')
plt.rcParams['figure.figsize'] = (12, 8)
```

```
visualizer = ClassBalance(labels=['Stayed', 'quit'])
visualizer.fit(hs.quit)
visualizer.show()
```

value 0 value 1  
kisi 0 or 1

**Apply** Stratified Sampling when creating training & validation set  
**Solution** → Apne example mein b bar graph mein bahut difference hai. So, we have a class imbalance problem. more employees who stayed than quit.

**Task 5 - Create Training & Test sets**

Separate data into feature matrix  $X$  and target vector  $y$ .  $X$  is called  $X$ .  
 size of quit ki hai target.

```
X = hs.loc[:, hs.columns != 'quit']
y = hs.quit
```

locate not quit

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(
    X, y, random_state=0, test_size=0.2, stratify=y)
```



stratify y ensures class distribution is approx that of y. i.e.  
 Eg - Y contains 25% - '0' and 75% - '1' then g-bran by test will have 25% - '0' and 75% - '1'

**TASK 6 → Build Decision Tree Classifier with Interactive console**

To make interactive we use @interact.

```
@interact
def plot-tree (crit=['gini', 'entropy'],
               splits=['best', 'random'],
               depth = IntSlider(min=1, max=30, value=2,
                                continuous_update=False)):
```

- Dropdown

Aur yhi  
check code  
karo

less computation  
and overhead.

False → jab drag krke release  
krge tab update hoga

True → jaise drag krge vaise  
update hoga

Scikit-learn decision

tree classifier → estimator = DecisionTreeClassifier (random-state=0,  
 criterion=crit,  
 splitter=splits,  
 max\_depth=depth,  
 min\_samples\_split=min\_split,

min\_samples\_leaf = min\_leaf)

estimator.fit(x\_train, y\_train) ← Training data pe  
 print ('Decision Tree Training Accuracy = %.3f' %  
 format(accuracy\_score(y\_train,  
 estimator.predict(x\_train))

∴ 3/ matlab upto 3 decimal places

accuracy score → sklearn mein hai

y\_train → because training accuracy

estimator.predict(x\_train) → kya kyunhi fit krliya  
 huchka.



<sup>package name</sup>  
graph = Source (tell.export - graphviz / estimator,  
<sup>helper function</sup>  
<sup>predicted model</sup>  
no output file → out\_file ~~out\_file~~ = None,  
feature\_names = X\_train.columns,  
class\_names = ['Stayed', 'Quit'],  
display values in keep → filled = True)

display (Image (data = graph.pipe (format = 'png')))

**TASK 8 → Build Interactive Random Forest Classifier**

See Code Comments.

**TASK-9 Feature Estimator**

and  
Task 6 & 8 mein estimator return krwaao

ef = Copy returned estimator from  
Random forest (Task-8)

**Visualize**  
viz = FeatureImportances (ef)  
viz.fit = (X\_train, y\_train)  
viz.show ();