## STEP 1: Load the Dataset



## STEP 2: Clean & Prepare the Dataset



Remove duplicates

No duplicate rows were found.

1470 unique rows remain.

OK

## STEP 3: Convert CSV to ARFF

( actually not needed )

## STEP 4: Load Dataset in Weka Explorer



## STEP 5: Explore Attributes & Visualize Patterns

○ Attrition vs Job Satisfaction

## Attrition vs OverTime



## Histograms

## STEP 6: Create Visual Narratives in Flourish



| Attrition | Department | COUNTA of Attrition |
|---|---|---|
| ▬ | | 0 |
| Total | | 0 |
| ▬ No | Human Resourc | 51 |
| | Research & Dev | 828 |
| | Sales | 354 |
| No Total | | 1233 |
| ▬ Yes | Human Resourc | 12 |
| | Research & Dev | 133 |
| | Sales | 92 |
| Yes Total | | 237 |
| Grand Total | | 1470 |





Click for the story : https://public.flourish.studio/story/3459871/

**Screen shots :**

This visualization shows the distribution of attrition across departments. The Research & Development department has the highest number of exiting employees, followed by Sales. Human Resources experiences the lowest attrition.

## Attrition per department



The pie chart illustrates the overall share of attrition from each department. Research & Development accounts for the majority of employee exits, followed by Sales, while HR contributes the least.

## Attrition per department

## Step 7: Train Baseline Model

```
Test mode:    evaluate on training data

=== Classifier model (full training set) ===

RandomForest

Bagging with 100 iterations and base learner

weka.classifiers.trees.RandomTree -K 0 -M 1.0 -V 0.001 -S 1 -do-not-check-capabilities

Time taken to build model: 2.98 seconds

=== Evaluation on training set ===

Time taken to test model on training data: 0.33 seconds
```

## === Summary ===

| | |
|---|---|
| Correlation coefficient | 0.97981 |
| Mean absolute error | 0.6278 |
| Root mean squared error | 0.8194 |
| Relative absolute error | 30.4805 % |
| Root relative squared error | 32.8113 % |
| Total Number of Instances | 1470 |

Number of Leaves  :      2225
Size of the tree :       2247

Time taken to build model: 0.05 seconds

## === Evaluation on training set ===

Time taken to test model on training data: 0.07 seconds

## === Summary ===

| | | |
|---|---|---|
| Correctly Classified Instances | 1075 | 72.9308 % |
| Incorrectly Classified Instances | 399 | 27.0692 % |
| Kappa statistic | 0.6861 | |
| Mean absolute error | 0.0377 | |
| Root mean squared error | 0.1373 | |
| Relative absolute error | 39.0965 % | |
| Root relative squared error | 62.5605 % | |
| Total Number of Instances | 1474 | |

```
=== Detailed Accuracy By Class ===

                TP Rate  FP Rate  Precision  Recall  F-Measure  MCC     ROC Area  PRC Area  Class
                0.810    0.014    0.926      0.810   0.864      0.840   0.943     0.852     0
                0.211    0.004    0.762      0.211   0.330      0.386   0.804     0.299     1
                0.852    0.087    0.749      0.852   0.797      0.733   0.957     0.836     2
                0.134    0.000    1.000      0.134   0.236      0.350   0.934     0.525     3
                0.929    0.095    0.410      0.929   0.569      0.580   0.958     0.555     4
                0.613    0.001    0.905      0.613   0.731      0.740   0.971     0.691     5
                0.414    0.002    0.800      0.414   0.545      0.570   0.961     0.512     6
                0.856    0.049    0.749      0.856   0.799      0.764   0.976     0.825     7
                0.944    0.048    0.608      0.944   0.740      0.736   0.989     0.852     8
                0.922    0.007    0.855      0.922   0.887      0.883   0.999     0.960     9
                0.556    0.001    0.882      0.556   0.682      0.696   0.997     0.832     10
                0.727    0.002    0.842      0.727   0.780      0.780   0.999     0.890     11
                0.611    0.001    0.917      0.611   0.733      0.746   0.999     0.867     12
                0.571    0.000    1.000      0.571   0.727      0.754   0.999     0.878     13
                0.400    0.000    1.000      0.400   0.571      0.632   0.999     0.778     14
                0.800    0.000    1.000      0.800   0.889      0.894   1.000     0.967     15
                1.000    0.000    1.000      1.000   1.000      1.000   1.000     1.000     16
                0.714    0.000    1.000      0.714   0.833      0.845   1.000     0.937     17
Weighted Avg.   0.729    0.040    0.794      0.729   0.706      0.698   0.955     0.761

=== Confusion Matrix ===

   a   b   c   d   e   f   g   h   i   j   k   l   m   n   o   p   q   r   <-- classified as
 213   0  20   0  15   0   0   9   6   0   0   0   0   0   0   0   0   0 |   a = 0
  13  16  28   0   7   0   0   8   4   0   0   0   0   0   0   0   0   0 |   b = 1
   0   0 293   0  42   0   0   5   4   0   0   0   0   0   0   0   0   0 |   c = 2
   0   0  50  19  65   0   0   5   3   0   0   0   0   0   0   0   0   0 |   d = 3
   2   0   0   0  91   0   0   2   3   0   0   0   0   0   0   0   0   0 |   e = 4
   0   1   0   0   0  19   0   7   4   0   0   0   0   0   0   0   0   0 |   f = 5
   0   0   0   0   0   0  12  12   5   0   0   0   0   0   0   0   0   0 |   g = 6
   1   1   0   0   0   1   0 185  26   1   1   0   0   0   0   0   0   0 |   h = 7
   0   1   0   0   0   1   0   4 101   0   0   0   0   0   0   0   0   0 |   i = 8
   0   0   0   0   1   0   1   2   1  59   0   0   0   0   0   0   0   0 |   j = 9
   1   0   0   0   1   0   1   2   6   1  15   0   0   0   0   0   0   0 |   k = 10
   0   2   0   0   0   0   0   0   2   2   0  16   0   0   0   0   0   0 |   l = 11
   0   0   0   0   0   0   0   1   1   3   1   1  11   0   0   0   0   0 |   m = 12
   0   0   0   0   0   0   0   2   0   1   0   2   1   8   0   0   0   0 |   n = 13
   0   0   0   0   0   0   0   1   0   2   0   0   0   0   2   0   0   0 |   o = 14
   0   0   0   0   0   0   0   1   0   0   0   0   0   0   0   4   0   0 |   p = 15
   0   0   0   0   0   0   0   0   0   0   0   0   0   0   0   0   6   0 |   q = 16
   0   0   0   0   0   0   1   1   0   0   0   0   0   0   0   0   0   5 |   r = 17
```

**Number of Leaves  :      2225**

**Size of the tree :         2247**

**Time taken to build model: 0.12 seconds**

# === Stratified cross-validation ===

# === Summary ===

Correctly Classified Instances          749             50.8141 %
Incorrectly Classified Instances       725             49.1859 %
Kappa statistic                        0.4318
Mean absolute error                    0.0609
Root mean squared error                0.1874
Relative absolute error                63.1553 %
Root relative squared error            85.3806 %
Total Number of Instances              1474

=== Detailed Accuracy By Class ===

| | TP Rate | FP Rate | Precision | Recall | F-Measure | MCC | ROC Area | PRC Area | Class |
|---|---|---|---|---|---|---|---|---|---|
| | 0.768 | 0.014 | 0.922 | 0.768 | 0.838 | 0.812 | 0.884 | 0.794 | 0 |
| | 0.000 | 0.002 | 0.000 | 0.000 | 0.000 | -0.011 | 0.591 | 0.064 | 1 |
| | 0.718 | 0.078 | 0.737 | 0.718 | 0.728 | 0.646 | 0.916 | 0.762 | 2 |
| | 0.268 | 0.083 | 0.257 | 0.268 | 0.262 | 0.182 | 0.867 | 0.339 | 3 |
| | 0.490 | 0.083 | 0.296 | 0.490 | 0.369 | 0.324 | 0.885 | 0.277 | 4 |
| | 0.000 | 0.003 | 0.000 | 0.000 | 0.000 | -0.008 | 0.786 | 0.075 | 5 |
| | 0.000 | 0.010 | 0.000 | 0.000 | 0.000 | -0.014 | 0.793 | 0.055 | 6 |
| | 0.630 | 0.092 | 0.540 | 0.630 | 0.581 | 0.505 | 0.880 | 0.485 | 7 |
| | 0.561 | 0.116 | 0.275 | 0.561 | 0.369 | 0.325 | 0.830 | 0.233 | 8 |
| | 0.031 | 0.021 | 0.063 | 0.031 | 0.042 | 0.014 | 0.784 | 0.128 | 9 |
| | 0.148 | 0.017 | 0.138 | 0.148 | 0.143 | 0.126 | 0.660 | 0.078 | 10 |
| | 0.091 | 0.012 | 0.100 | 0.091 | 0.095 | 0.082 | 0.664 | 0.061 | 11 |
| | 0.000 | 0.005 | 0.000 | 0.000 | 0.000 | -0.008 | 0.592 | 0.029 | 12 |
| | 0.143 | 0.007 | 0.167 | 0.143 | 0.154 | 0.147 | 0.660 | 0.060 | 13 |
| | 0.000 | 0.003 | 0.000 | 0.000 | 0.000 | -0.003 | 0.591 | 0.015 | 14 |
| | 0.000 | 0.003 | 0.000 | 0.000 | 0.000 | -0.003 | 0.590 | 0.019 | 15 |
| | 0.833 | 0.001 | 0.833 | 0.833 | 0.833 | 0.833 | 0.999 | 0.924 | 16 |
| | 0.429 | 0.001 | 0.600 | 0.429 | 0.500 | 0.505 | 0.779 | 0.282 | 17 |
| Weighted Avg. | 0.508 | 0.058 | 0.495 | 0.508 | 0.494 | 0.443 | 0.847 | 0.479 | |

=== Confusion Matrix ===

| a | b | c | d | e | f | g | h | i | j | k | l | m | n | o | p | q | r | <-- classified as |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 202 | 0 | 18 | 7 | 15 | 0 | 0 | 11 | 9 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | a = 0 |
| 13 | 0 | 22 | 14 | 5 | 1 | 0 | 11 | 9 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | b = 1 |
| 0 | 0 | 247 | 48 | 32 | 0 | 0 | 5 | 11 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | c = 2 |
| 0 | 0 | 36 | 38 | 57 | 0 | 0 | 6 | 4 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | d = 3 |
| 1 | 0 | 3 | 33 | 48 | 0 | 0 | 3 | 9 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | e = 4 |
| 0 | 1 | 2 | 7 | 3 | 0 | 0 | 8 | 8 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | f = 5 |
| 0 | 0 | 1 | 0 | 0 | 0 | 0 | 12 | 11 | 0 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | g = 6 |
| 1 | 1 | 0 | 0 | 0 | 1 | 0 | 136 | 56 | 8 | 6 | 4 | 0 | 1 | 1 | 1 | 0 | 0 | h = 7 |
| 0 | 0 | 3 | 0 | 0 | 0 | 3 | 21 | 60 | 9 | 5 | 3 | 1 | 0 | 0 | 2 | 0 | 0 | i = 8 |
| 1 | 0 | 0 | 0 | 0 | 0 | 5 | 19 | 28 | 2 | 2 | 3 | 1 | 1 | 1 | 0 | 1 | 0 | j = 9 |
| 0 | 0 | 0 | 0 | 1 | 0 | 3 | 6 | 4 | 4 | 4 | 1 | 1 | 1 | 1 | 0 | 0 | 1 | k = 10 |
| 0 | 1 | 1 | 0 | 0 | 2 | 0 | 5 | 3 | 3 | 1 | 2 | 2 | 2 | 0 | 0 | 0 | 0 | l = 11 |
| 1 | 0 | 0 | 0 | 0 | 0 | 0 | 4 | 1 | 4 | 2 | 4 | 0 | 2 | 0 | 0 | 0 | 0 | m = 12 |
| 0 | 0 | 1 | 0 | 0 | 0 | 1 | 2 | 1 | 0 | 3 | 2 | 2 | 2 | 0 | 0 | 0 | 0 | n = 13 |
| 0 | 0 | 0 | 1 | 0 | 0 | 0 | 2 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | o = 14 |
| 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 2 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | p = 15 |
| 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 5 | 0 | q = 16 |
| 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 3 | r = 17 |

# Overfitting

## Using naïve bayes as training results

```
=== Evaluation on training set ===

Time taken to test model on training data: 0.3 seconds

=== Summary ===

Correctly Classified Instances        1258               85.346 %
Incorrectly Classified Instances       216               14.654 %
Kappa statistic                          0.8313
Mean absolute error                      0.0214
Root mean squared error                  0.1089
Relative absolute error                 22.2378 %
Root relative squared error             49.6316 %
Total Number of Instances             1474
```

## Using naïve bayes for cross validation

```
=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances         742               50.3392 %
Incorrectly Classified Instances       732               49.6608 %
Kappa statistic                          0.4244
Mean absolute error                      0.0586
Root mean squared error                  0.1992
Relative absolute error                 60.7984 %
Root relative squared error             90.7469 %
Total Number of Instances             1474
```

## By using random forest without converting all data to nominal and removing unrelated rows

## === Detailed Accuracy By Class ===

| TP Rate | FP Rate | Precision | Recall | F-Measure | MCC | ROC Area | PRC Area | Class |
|---------|---------|-----------|--------|-----------|--------|----------|----------|-------|
| 0.251 | 0.058 | 0.485 | 0.251 | 0.331 | 0.255 | 0.629 | 0.299 | 0 |
| 0.013 | 0.014 | 0.048 | 0.013 | 0.021 | -0.002 | 0.449 | 0.048 | 1 |
| 0.767 | 0.630 | 0.271 | 0.767 | 0.401 | 0.123 | 0.614 | 0.310 | 2 |
| 0.063 | 0.036 | 0.158 | 0.063 | 0.090 | 0.042 | 0.572 | 0.132 | 3 |
| 0.102 | 0.019 | 0.278 | 0.102 | 0.149 | 0.134 | 0.595 | 0.109 | 4 |
| 0.032 | 0.007 | 0.091 | 0.032 | 0.048 | 0.042 | 0.544 | 0.027 | 5 |
| 0.034 | 0.005 | 0.125 | 0.034 | 0.054 | 0.056 | 0.579 | 0.030 | 6 |
| 0.199 | 0.077 | 0.307 | 0.199 | 0.242 | 0.147 | 0.622 | 0.215 | 7 |
| 0.028 | 0.023 | 0.086 | 0.028 | 0.042 | 0.008 | 0.571 | 0.085 | 8 |
| 0.047 | 0.011 | 0.158 | 0.047 | 0.072 | 0.064 | 0.629 | 0.077 | 9 |
| 0.037 | 0.008 | 0.077 | 0.037 | 0.050 | 0.041 | 0.616 | 0.053 | 10 |
| 0.000 | 0.005 | 0.000 | 0.000 | 0.000 | -0.009 | 0.588 | 0.023 | 11 |
| 0.000 | 0.002 | 0.000 | 0.000 | 0.000 | -0.005 | 0.531 | 0.024 | 12 |
| 0.071 | 0.003 | 0.200 | 0.071 | 0.105 | 0.115 | 0.628 | 0.046 | 13 |
| 0.000 | 0.001 | 0.000 | 0.000 | 0.000 | -0.002 | 0.283 | 0.003 | 14 |
| 0.000 | 0.002 | 0.000 | 0.000 | 0.000 | -0.003 | 0.565 | 0.005 | 15 |
| 0.000 | 0.001 | 0.000 | 0.000 | 0.000 | -0.001 | 0.228 | 0.001 | 16 |
| 0.000 | 0.001 | 0.000 | 0.000 | 0.000 | -0.002 | 0.689 | 0.039 | 17 |
| **Weighted Avg.** | **0.274** | **0.177** | **0.252** | **0.274** | **0.218** | **0.116** | **0.597** | **0.193** |

## === Summary ===

| | | |
|---|---|---|
| **Correctly Classified Instances** | 403 | 27.415 % |
| **Incorrectly Classified Instances** | 1067 | 72.585 % |
| **Kappa statistic** | 0.099 | |
| **Mean absolute error** | 0.089 | |
| **Root mean squared error** | 0.231 | |
| **Relative absolute error** | 92.3523 % | |
| **Root relative squared error** | 105.2777 % | |
| **Total Number of Instances** | 1470 | |

# Underfitting

**Using zeroR classifiervto test results**

**As a training set**

```
=== Evaluation on training set ===

Time taken to test model on training data: 0 seconds

=== Summary ===

Correctly Classified Instances         344              23.3379 %
Incorrectly Classified Instances      1130              76.6621 %
Kappa statistic                          0
Mean absolute error                      0.0964
Root mean squared error                  0.2195
Relative absolute error                100      %
Root relative squared error            100      %
Total Number of Instances             1474
```

## Cross validating

```
=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances         344              23.3379 %
Incorrectly Classified Instances      1130              76.6621 %
Kappa statistic                          0
Mean absolute error                      0.0965
Root mean squared error                  0.2195
Relative absolute error                100      %
Root relative squared error            100      %
Total Number of Instances             1474
```

---

## Step 10: Summarize Insights : using napkin

Napkin ai summary link :

**https://app.napkin.ai/page/CgoiCHByb2Qtb25lEiwKBFBhZ2UaJDc5ZDE5MDc3LWMxMjUtNGEyMi1hNWZhLWY3NDQ2ODEyMmVhNA?s=1**

summary main points

- The evaluation of models (J48, Naive Bayes, Random Forest, ZeroR) showed mixed performance, with J48 suffering from **overfitting** and Random Forest and ZeroR showing **underfitting**, indicating improper learning of meaningful patterns.

- A major challenge identified across models was **class imbalance**, where "no attrition" cases dominate the dataset, causing models to bias predictions toward the majority class and reduce detection accuracy of actual attrition cases.
- The inclusion of irrelevant features like **EmployeeNumber** negatively impacted learning by creating meaningless model splits; removing such non-predictive features is essential to improve model interpretability and accuracy.
- The summary highlighted the need for techniques like **resampling, cost-sensitive learning, proper feature selection, and hyperparameter tuning** to improve prediction performance and handle imbalanced data effectively.
- Potential key drivers of attrition include factors such as **job satisfaction, salary, work-life balance, tenure, promotion history, and workload (overtime)**, although further feature importance analysis is needed once the model is improved.
- Based on insights, HR can leverage refined predictive models to identify high-risk employees and implement targeted retention strategies (e.g., salary review, career development programs, flexible working policies, and department-specific interventions).

---

## Reflection Questions

Answer any 3:

1. Which factors contributed most to attrition?
2. Did your model overfit or underfit? Why?
3. How can dataset quality be improved?
4. How should HR teams use this model?

Ans 1 - From the patterns observed, employees with low job satisfaction, poor work-life balance, lower salaries, and limited growth opportunities seemed more likely to leave. Tenure and overtime also appeared to influence the likelihood of attrition.

Ans2 - Yes — one model overfitted and another underfitted. The J48 model overfitted because it learned the training data too perfectly, including noise. Meanwhile, Random Forest and ZeroR underfitted due to class imbalance and irrelevant features affecting learning.

Ans3 - The dataset can be improved by removing non-useful columns like EmployeeNumber, balancing the classes, and cleaning or encoding categorical values properly. Doing this will help the model learn real patterns instead of random or misleading ones.

---

Overall, the mini-project demonstrated how different machine learning models behave when applied to HR attrition data. The experiment helped me understand how data imbalance, irrelevant features, and model settings affect prediction accuracy. After improving feature selection and handling imbalance, the model can become more useful for HR teams to identify employees at risk and take preventive actions.

---

**Final dataset :**

---

**Submitted by** <u>Abduttaiyeb Huseni Matcheswala</u> **(b24bs1015)**