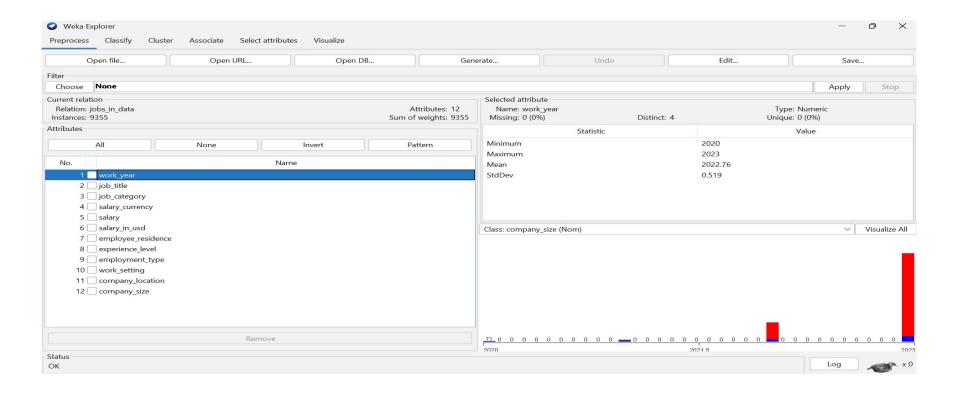# HW4 - Classification

Team Members: Puja Shah & Sanjida Chowdhury
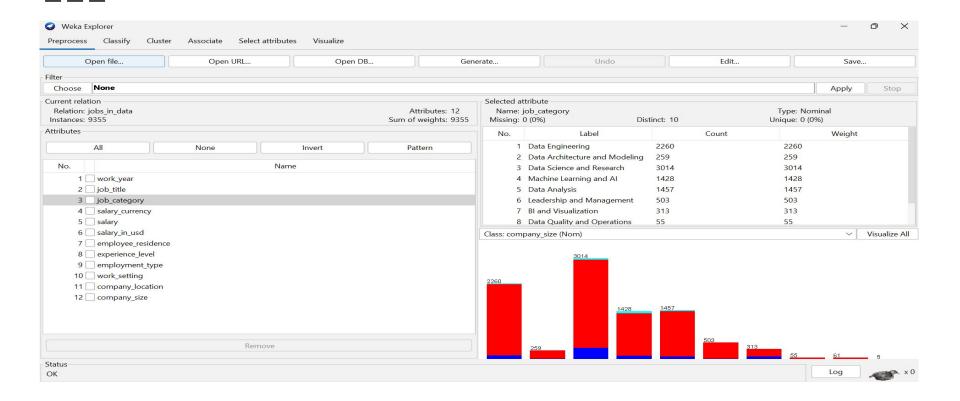
# Dataset: "Jobs and Salaries in Data Science"

---

# Dataset: "Jobs and Salaries in Data Science"

_ _ _

- The dataset used for this analysis contains information related to job postings, including attributes such as work year, job title, job category, salary, experience level, and more. It comprises a total of 9355 instances, each representing a unique job posting. The dataset attributes provide valuable insights into the characteristics of different job opportunities, such as the type of role, salary range, and required experience level. This information can be valuable for job seekers, employers, and researchers interested in understanding trends and patterns in the job market.
- The "Jobs and Salaries in Data Science" dataset sourced from Kaggle (same data from the pre-processing). The link to the data set is here: https://www.kaggle.com/datasets/hummaamqaasim/jobs-in-data
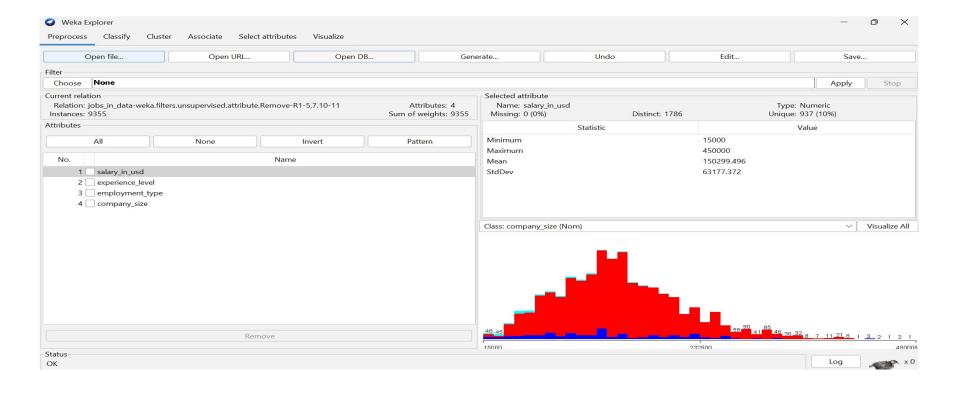
# Distribution of the Target Class

– – –

# Distribution of the Target Class

———

- The target class for this classification task is the job category.
- It serves as the class label that we aim to predict or classify.
- The job category is chosen because the goal is to categorize job listings into different job categories.
- Attributes like salary, experience level, employment type, and company size are used as features for classification.
- These attributes help determine the nature of the job and are input features for the classification model.
- Therefore, the job category becomes the target class as it encapsulates the main objective of classifying jobs into specific categories.
- By predicting the job category, we can effectively categorize job listings based on various attributes.
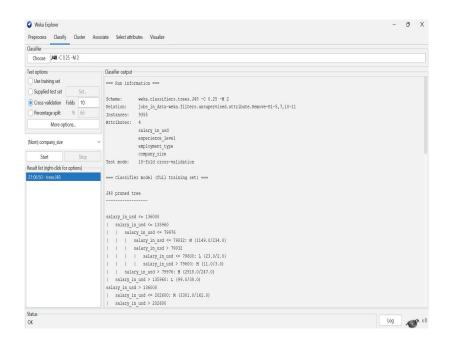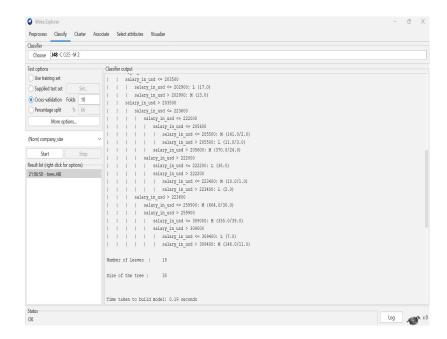
# Attributes to Train the Model

– – –

# Attributes to Train the Model

— — —

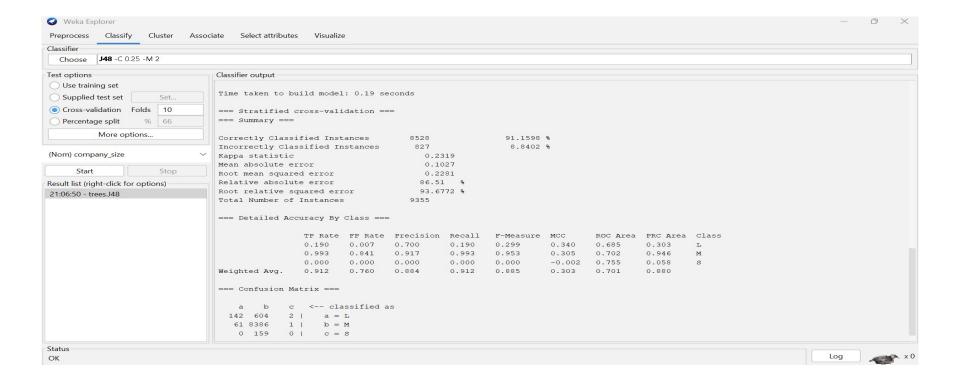Chosen attributes for training the model are:

- **salary_in_usd:** Represents standardized salary in USD, crucial for job seekers and employers.
- **experience_level:** Indicates the level of experience required, impacting salary and job responsibilities.
- **employment_type:** Specifies employment terms (e.g., Full-time, Part-time), influencing work conditions.
- **company_size:** Denotes the size of the company, providing insights into organizational dynamics and career growth opportunities.
- These attributes were selected because they collectively cover essential aspects of job listings relevant to both job seekers and employers, including financial factors, job requirements, employment terms, and organizational characteristics. Other attributes excluded because of less informative for predicting job suitability and salary expectations.

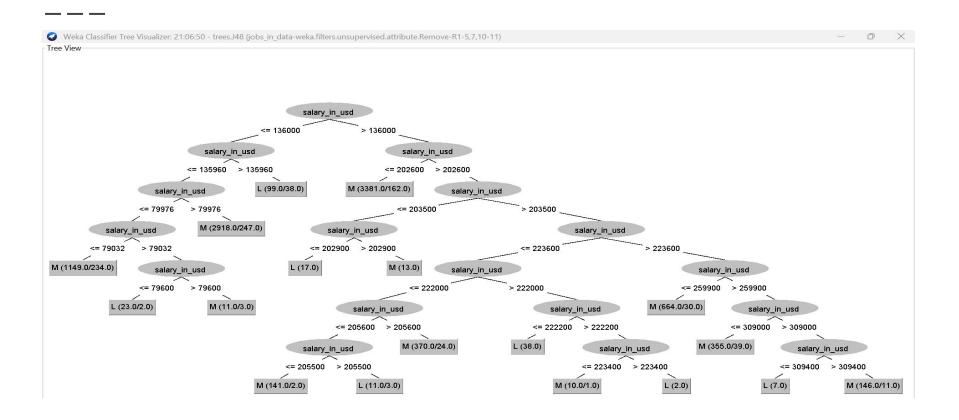# Classifier Model: J48 -C 0.25 -M 2

# Classifier Model: J48 -C 0.25 -M 2

‒ ‒ ‒

# Decision Tree

# Interpreting a J48 Decision Tree Model

— — —

- The J48 decision tree model achieved an overall accuracy of approximately 91.16% in the stratified 10-fold cross-validation.
- The model's decision tree is relatively complex, with 35 nodes and 18 leaves, indicating multiple decision paths based on the attributes salary_in_usd, experience_level, employment_type, and company_size.
- Specific rules generated by the tree involve conditions on salary_in_usd to predict job categories. For example, if salary_in_usd is less than or equal to $136,000, the model predicts 'L' (low salary job). If salary_in_usd is greater than $136,000, the model further examines the salary range to make predictions.
- The model has 827 incorrect classifications out of 9355 instances in total, indicating that it misclassified approximately 8.84% of instances.
- The specific errors for each possible classification outcome vary: For class 'L' (low salary), there were 142 instances misclassified as 'L', 604 instances correctly classified as 'L', and 2 instances misclassified as 'M'.
- For class 'M' (medium salary), there were 61 instances misclassified as 'L', 8386 instances correctly classified as 'M', and 1 instance misclassified as 'S'.
- For class 'S' (low frequency), there were no instances correctly classified, indicating a significant error for this class.

# Evaluation and Summary

_ _ _

- Recall for class 'L' (low salary): 0.190
- Recall for class 'M' (medium salary): 0.993
- Recall for class 'S' (low frequency): 0.000
- Precision for class 'L': 0.700
- Precision for class 'M': 0.917
- Precision for class 'S': 0.000
- Accuracy: 91.16%
- The model performed relatively well in predicting classes with medium salaries (Class M), achieving a high recall and precision. However, it struggled with classes of low salaries (Class L) and low frequency (Class S), as indicated by low recall and precision scores.
- The model met expectations in some aspects, such as accurately predicting medium salary jobs, but it fell short in predicting low salary and low frequency jobs.
- The limited predictive power for low salary and low frequency jobs could be due to the complexity of real-world job market dynamics.

Thank You