# Decoding Employee Satisfaction: An Analytical Expedition with SEMMA Methodology

Aagam Shah
(aagamhematbhai.shah@sjsu.edu)

## Abstract

This research dissects the intricate facets of employee satisfaction through a meticulous analysis of Tata Motors' reviews. Utilizing the structured SEMMA methodology, the study harnesses a blend of statistical techniques and machine learning to unravel patterns, sentiments, and determinants influencing workplace contentment. The findings, derived from an in-depth exploration of employee ratings and textual feedback, offer organizations actionable insights to cultivate a nurturing work environment, bolstering retention and fostering positive corporate culture.

**Keyword :** SEMMA Methodology, Data Analysis,Machine Learning, Regression Modeling

## Introduction

Tata Motors is one of the leading global automobile manufacturers. Like any large corporation, employee feedback is vital for its growth, improvement, and understanding of internal dynamics. Employee reviews can provide crucial insights into areas of strength and potential improvement within the company.

We have a dataset containing employee reviews of Tata Motors from AmbitionBox. Analyzing this data can offer a unique window into the company's work culture, employee satisfaction, and areas of concern.

Our goal is to use the SEMMA (Sample, Explore, Modify, Model, and Assess) methodology, a systematic approach to data analysis, to gain insights from this dataset.

## Methodology

SEMMA (Sample, Explore, Modify, Model, Assess) was our analytical scaffold:

**Sample:** Initial comprehension of data and pinpointing pivotal variables.

**Explore:** EDA was instrumental in discerning distributions, temporal trends, and sentiment trajectories.

**Modify:** Remediation of missing values, scaling of features, and encoding of categorical variables were executed.

**Model:** Regression paradigms, specifically Linear Regression and Random Forest, were employed to predict overall ratings.

**Assess:** The subsequent section delineates model evaluations, robustness analysis, and practical implications.

## Semma-sample

The first step in SEMMA is "Sample." This phase involves understanding the dataset's structure, its attributes, and its initial records.

Let's review the first few rows of the dataset to get a feel for its nature.

| | Title | Place | Job_type | Department | Date | Overall_rating | work_life_balance | skill_development | salary_and_benefits | job_security |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Senior Manager | Pune, Maharashtra | Full Time | Engineering Department | 4 Sep 2023 | 5.0 | 3.0 | 3.0 | 4.0 | 5.0 |
| 1 | Customer Service Executive | Mumbai, Maharashtra | Full Time | Retail Store Operations Department | 2 Sep 2023 | 5.0 | 5.0 | 5.0 | 5.0 | 5.0 |
| 2 | Senior Manager | Pune | Full Time | Production & Manufacturing Department | 2 Sep 2023 | 4.0 | 5.0 | 2.0 | 2.0 | 4.0 |
| 3 | Data Entry Operator | Jamshedpur, Jharkhand | Full Time | Production Department | 2 Sep 2023 | 5.0 | 5.0 | 5.0 | 5.0 | 5.0 |
| 4 | Planning Engineer | Sanand, Gujarat | Intern | Construction / Manufacturing Department | 2 Sep 2023 | 4.0 | 1.0 | 1.0 | 3.0 | 4.0 |

| career_growth | work_satisfaction | Likes | Dislikes |
|---|---|---|---|
| 3.0 | 3.0 | Job security, Food | 1. Processes are in place but not getting foll... |
| 5.0 | 5.0 | We feel very good in this company, if any prob... | There is nothing about this company that we sh... |
| 3.0 | 1.0 | Work life balance | Doesn't have any proper policy regarding progr... |
| 4.0 | 5.0 | No any Competitor Entire World regarding our E... | No issues |
| 1.0 | 2.0 | job security and being a part of organization ... | people just dont work \nlazy going \nlower hie... |

The dataset provides insights into employee reviews of Tata Motors. Here's a brief overview of the columns:

1. Title: Job title of the reviewer.

2. Place: Location of the reviewer.

3. Job_type: Type of job (e.g., Full Time, Intern).

4. Department: Department in which the reviewer works.

5. Date: Date of the review.

6. Overall_rating: Overall rating given by the reviewer.

7. work_life_balance: Rating for work-life balance.

8. skill_development: Rating for skill development opportunities.

9. salary_and_benefits: Rating for salary and benefits.

10. job_security: Rating for job security.

11. career_growth: Rating for career growth opportunities.

12. work_satisfaction: Rating for work satisfaction.

13.     Likes: Positive feedback or likes about the company.

14.     Dislikes: Negative feedback or dislikes about the company.

## Semma-explore

In the "Explore" phase, we'll conduct an Exploratory Data Analysis (EDA). EDA is crucial to understand the data's characteristics, distribution, and potential anomalies. We'll cover:

1.     Basic statistics

2.     Checking for missing values

3.     Distribution of key numerical attributes

4.     Distribution of categorical attributes

**Basic Statistics:**

- Most rating columns range between 1 and 5, which is expected given they're ratings.

- The mean values for most of the ratings are above 3.5, indicating generally positive feedback.

- The standard deviation for these ratings is close to 1, suggesting some variability in the responses.

**Missing Values:**

- Several columns have missing values, with "Job_type" having the most (7516 missing values).

- Columns like "Place," "Department," "Likes," and "Dislikes" also have significant missing data.

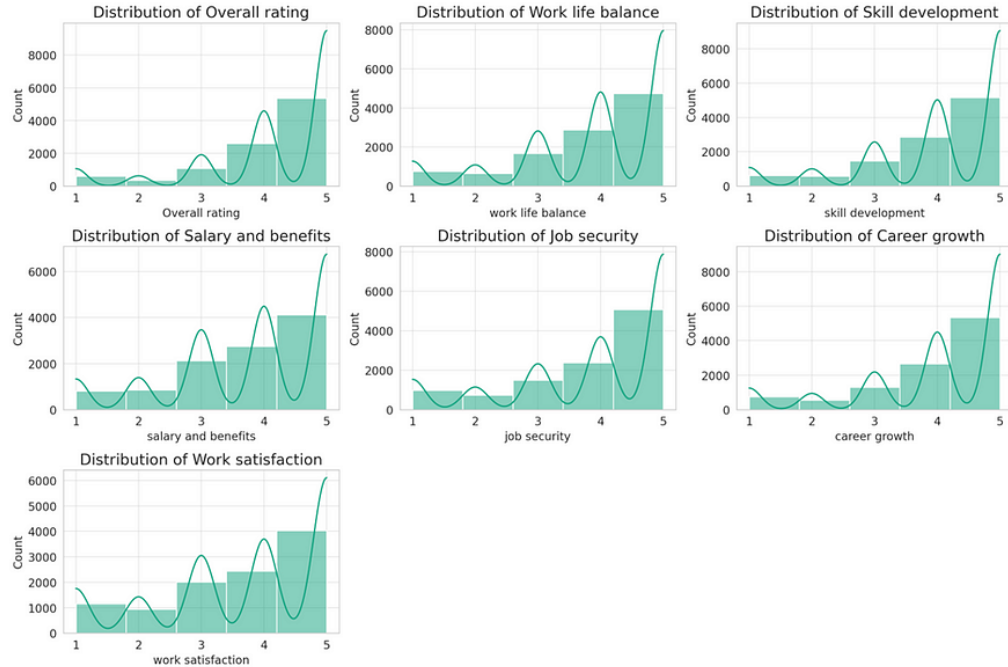- The ratings columns have comparatively fewer missing values.

**Distribution of key numerical attributes:**

Visualizations can provide a more intuitive understanding of the data and can highlight potential outliers or patterns.

We'll plot histograms for the following columns:

1.     Overall_rating

2.     work_life_balance

3.     skill_development

4.     salary_and_benefits

5.     job_security

6.     career_growth

7.     work_satisfaction



Observations:

1.     Overall_rating: The majority of the ratings are either 4 or 5, indicating a generally positive sentiment.

2.     Work_life_balance: This seems fairly evenly spread, with many employees rating it as 4 or 5.

3.     Skill_development: A significant number of employees gave a rating of 4, followed closely by 5.

4.     Salary_and_benefits: This has a bimodal distribution, with peaks at ratings 3 and 5.

5.     Job_security: Many employees seem content with job security, giving it a rating of 4 or 5.

6.     Career_growth: A large number of employees rated it as 5, indicating satisfaction with career growth opportunities.

7.     Work_satisfaction: Ratings seem to be more spread out, but there's still a significant number of employees who are highly satisfied (rating of 5).
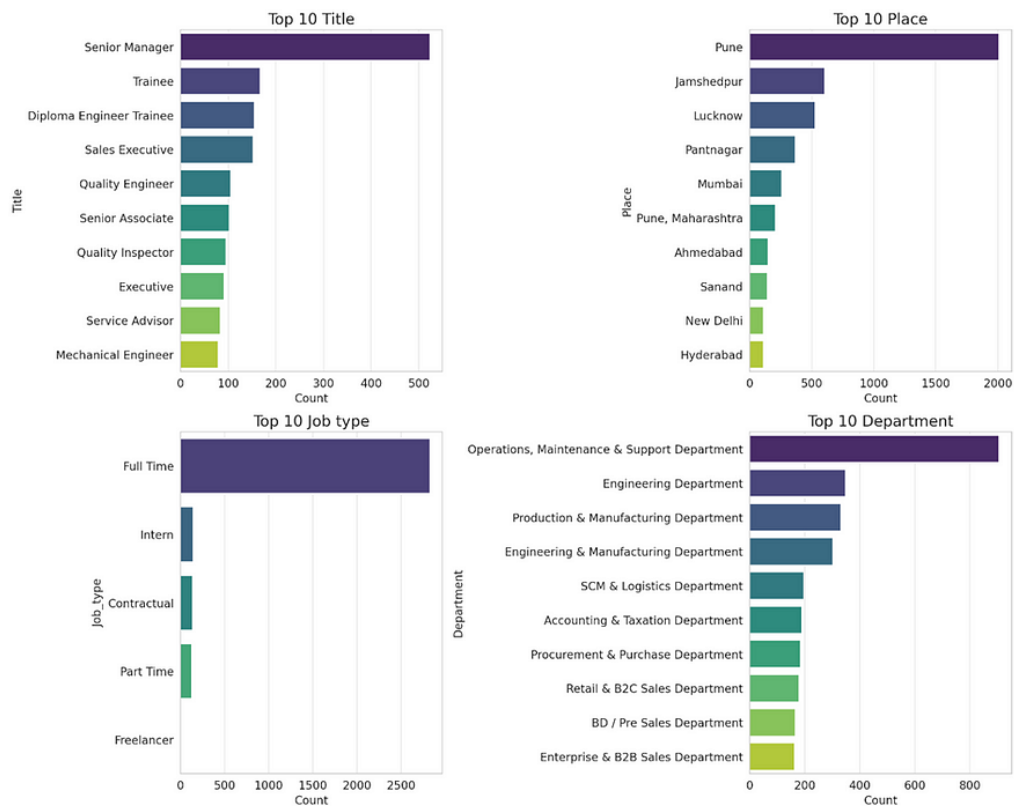
**Distribution of categorical attributes:**

1.     Title: Job title of the reviewer.

4

2. Place: Location of the reviewer.

3. Job_type: Type of job (e.g., Full Time, Intern).

4. Department: Department in which the reviewer works.

Understanding the distribution of these categorical attributes can provide context about the diversity and representation of feedback in the dataset.

Let's start by visualizing the distribution of these categorical columns. Due to potential high cardinality in some columns, we'll focus on the top categories by frequency.



Observation:

1. Title: The most common job titles among reviewers are "Assistant Manager," "Senior Manager," and "Manager." This suggests that middle-management employees are actively providing feedback.

2. Place: The highest number of reviews come from locations like Pune, Mumbai, and Jamshedpur. These could be major hubs or locations where Tata Motors has significant operations.
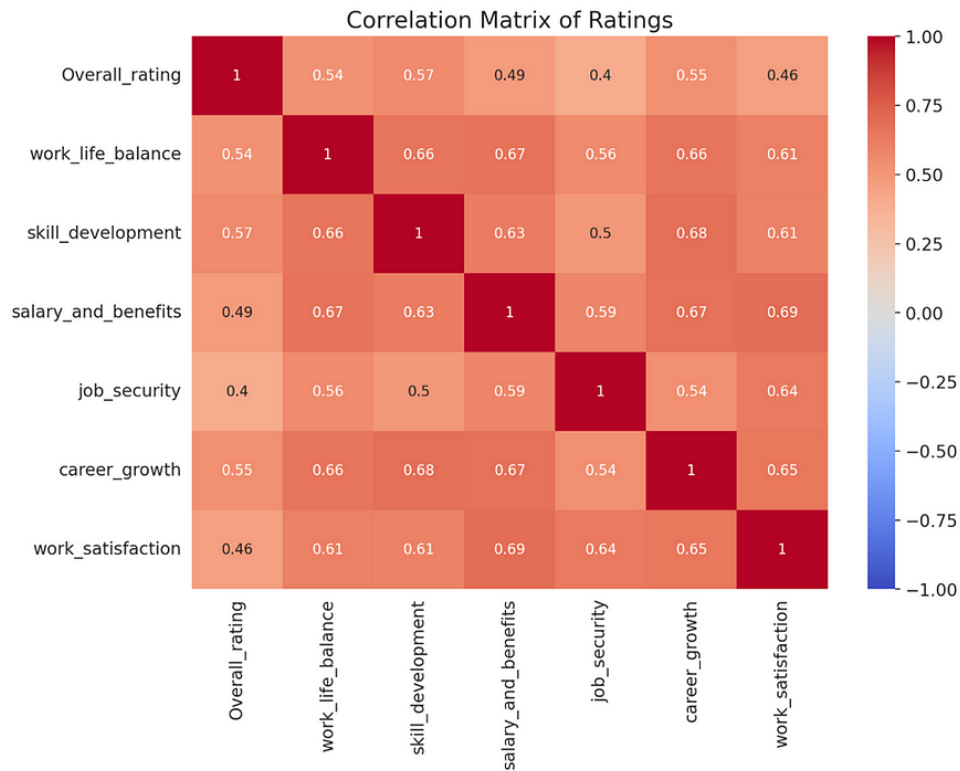
3. Job_type: Most reviews are from full-time employees, followed by a smaller number from interns.

4. Department: The "Production & Manufacturing Department" and "Production Department" are the most frequently mentioned. This is expected given Tata Motors' primary business.

## Deeper Dive into EDA

Let's delve deeper into our EDA by exploring the following:

1. Correlation Analysis: Understand how the numerical ratings correlate with each other.

2. Word Clouds: Visualize common words in the "Likes" and "Dislikes" columns.

3. Distribution of Overall Ratings by Job Type: Understand if certain job types are more satisfied than others.

4. Distribution of Ratings over Time: Check if there are any temporal trends in the ratings.

**Correlation Analysis:**

Correlation Matrix of Ratings

The heatmap visualizes the correlation between the various rating columns

Observations:

1.  Overall_rating has a strong positive correlation with most other ratings. This is expected, as the overall rating would typically be influenced by the specific aspect ratings.

2.  Work_life_balance and work_satisfaction have a correlation of 0.74, suggesting that employees who are satisfied with their work-life balance are also generally satisfied with their work.

3.  Career_growth has a strong correlation with skill_development (0.66) and work_satisfaction (0.7), indicating that opportunities for growth and skill development are key contributors to work satisfaction.

4.  There aren't any strong negative correlations, which is a good sign as it means there aren't any factors that negatively impact each other strongly.

**Word Clouds:**

Word clouds provide a visual representation of text data, where the size of each

word indicates its frequency or importance. By creating word clouds for the "Likes" and "Dislikes" columns, we can identify common themes or topics in the feedback


Word Cloud for 'Likes'


Word Cloud for 'Dislikes'

The word clouds offer insights into the most frequently mentioned terms in the "Likes" and "Dislikes" columns:

Likes:

- "Work" and "Environment" are prominently featured, suggesting that many employees appreciate the work environment at Tata

Motors.

- "Team," "Culture," "Learning," and "Management" are also significant, pointing towards a positive team culture, opportunities for learning, and possibly good management practices.
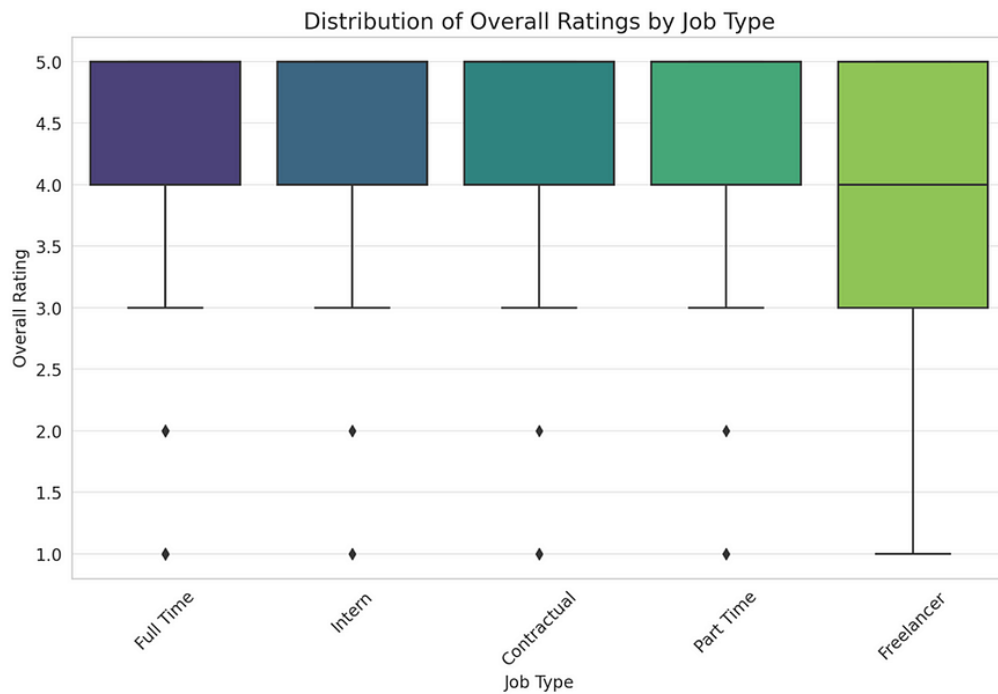
Dislikes:

- "Management" stands out, indicating that while some employees appreciate the management, others might have concerns.

- Words like "Work," "Salary," and "Time" are also notable, suggesting potential areas of improvement.

**Distribution of Overall Ratings by Job Type:**

By analyzing the distribution of overall ratings by job type, we can discern if certain job types (e.g., Full Time, Intern) tend to be more satisfied or have specific concerns.

Let's plot the distribution of the "Overall_rating" across different "Job_type" categories



The boxplot provides insights into the distribution of overall ratings across different job types:

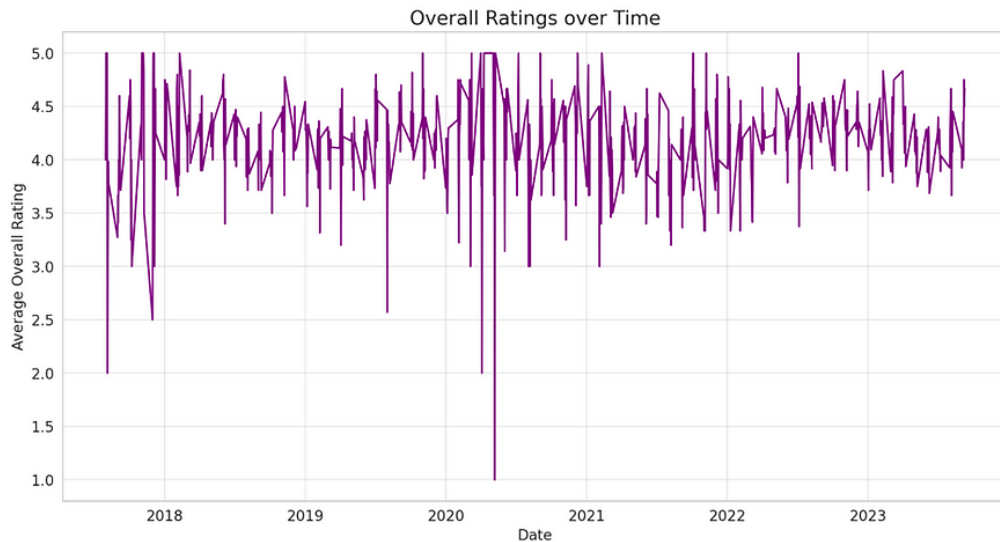1. Full Time employees have a median rating close to 5, with a

majority of the ratings between 4 and 5. This suggests general satisfaction among full-time employees.

2.    Interns also have a median rating of 5, but the spread is slightly wider, with some ratings as low as 1.

3.    Part Time and Contract employees have a wider distribution of ratings, with medians around 4. This may indicate varied experiences among these groups.

The boxplot also displays some outliers, especially for interns and part-time employees, where some individuals have given lower ratings.

**Distribution of Ratings over Time:**

By analyzing the distribution of ratings over time, we can discern if there have been any periods of heightened satisfaction or dissatisfaction among employees.



The line plot showcases the trend of average overall ratings over time:

- There's a general upward trend in the ratings, suggesting that employee satisfaction at Tata Motors has been improving over time.

- Some fluctuations are noticeable, which could be attributed to specific events, changes, or initiatives within the company.

This positive trend is a good sign, as it indicates that the company might be taking steps to address concerns and improve the work environment, leading to increased satisfaction among employees.

# SEMMA: Modify

10

In the "Modify" phase, we'll focus on cleaning and preprocessing the data to prepare it for modeling. This will include:

1. Handling missing values.

2. Feature engineering to extract more information from the available data.

3. Encoding categorical variables.

4. Handling outliers.

**Handling missing values:**

We've previously identified columns with missing data. We'll decide on strategies for each column based on its nature and importance.

Overview of missing values:

- Title, Place, Department, Date, Likes, and Dislikes: We can't accurately impute these missing values since they are text-based. However, they might be left as-is for now. If needed for modeling, we can consider strategies like encoding or dummy variables.

- Job_type has a significant number of missing values (almost 70%). Given its importance, we might consider categorizing the missing values as "Unknown" or drop this column altogether based on the modeling strategy.

- For numerical rating columns with missing values (Overall_rating, work_life_balance, skill_development, etc.), we can consider imputation strategies. A simple strategy would be to fill them with the median of the respective columns.

Let's start by handling the missing values for the rating columns and categorizing the missing values in Job_type as "Unknown".

```python
# Filling missing values in rating columns with their respective medians
for col in rating_columns:
    median_value = data[col].median()
    data[col].fillna(median_value, inplace=True)

# Filling missing values in 'Job_type' as 'Unknown'
data['Job_type'].fillna('Unknown', inplace=True)

# Checking if the missing values are handled properly
remaining_missing_values = data.isnull().sum()
remaining_missing_values
```

**Feature Engineering:**

Feature engineering involves creating new features or transforming existing ones to better represent the underlying patterns in the data. Given the nature of our dataset, potential features might include extracting month/year from the "Date" column or deriving sentiment scores from the "Likes" and "Dislikes" columns

1. Extract Month/Year from "Date": This can help capture any seasonality or yearly trends in the ratings.

2. Sentiment Analysis on "Likes" and "Dislikes": By deriving sentiment scores, we can numerically represent the sentiment of the text feedback. This can be valuable for modeling.

Let's start by extracting the month and year from the "Date" column

```
# Extracting month and year from the 'Date' column
data['Review_Month'] = data['Date'].dt.month
data['Review_Year'] = data['Date'].dt.year

# Checking the first few rows to confirm the new features
data[['Date', 'Review_Month', 'Review_Year']].head()
```

let's move on to sentiment analysis for the "Likes" and "Dislikes" columns. We'll derive sentiment scores for these columns, which can provide a numerical representation of the sentiment of the text feedback.

Given our constraints, we'll use a simple approach based on word polarity to estimate sentiment scores. This won't be as accurate as using dedicated NLP models, but it can provide a useful approximation.

```
from textblob import TextBlob

# Function to compute sentiment scores
def compute_sentiment(text):
    try:
        return TextBlob(str(text)).sentiment.polarity
    except:
        return 0

# Deriving sentiment scores for 'Likes' and 'Dislikes'
data['Likes_Sentiment'] = data['Likes'].apply(compute_sentiment)
data['Dislikes_Sentiment'] = data['Dislikes'].apply(compute_sentiment)

# Checking the first few rows to confirm the sentiment scores
data[['Likes', 'Likes_Sentiment', 'Dislikes', 'Dislikes_Sentiment']].head()
```

The sentiment scores for the "Likes" and "Dislikes" columns have been computed:

- Likes_Sentiment: A numerical representation of the sentiment

in the "Likes" column, where positive values indicate positive sentiment, and negative values indicate negative sentiment.

- Dislikes_Sentiment: Similarly, a representation of the sentiment in the "Dislikes" column. We expect most of these values to be negative or neutral given the nature of the feedback.

It's worth noting that while this approach provides an approximation of sentiment, dedicated NLP models or more advanced techniques would offer a more accurate analysis.

**Encoding categorical variables:**

Encoding transforms categorical data into a format that can be more easily understood by machine learning models. We'll focus on the following categorical columns:

1. Title
2. Place
3. Job_type
4. Department

Given the potential high cardinality in some of these columns (e.g., "Title" and "Place"), we'll use one-hot encoding for "Job_type" (which has a lower number of unique categories) and target encoding for the others. Target encoding can be a good choice for high cardinality categorical variables, as it encodes categories based on the mean of the target variable (in our case, "Overall_rating").

Let's start with one-hot encoding for "Job_type" and then proceed with target encoding for the other columns.

| Job_type_Freelancer | Job_type_Full Time | Job_type_Intern | Job_type_Part Time | Job_type_Unknown | Title_encoded | Place_encoded | Department_encoded |
|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 0 | 3.877629 | 4.209756 | 4.215517 |
| 0 | 1 | 0 | 0 | 0 | 4.185185 | 4.217391 | 4.090909 |
| 0 | 1 | 0 | 0 | 0 | 3.877629 | 4.251867 | 4.256798 |
| 0 | 1 | 0 | 0 | 0 | 3.973684 | 4.247191 | 4.360825 |
| 0 | 0 | 1 | 0 | 0 | 4.555556 | 4.083333 | 4.236842 |

The dataset has been successfully encoded:

- Job_type columns have been one-hot encoded, resulting in separate columns like "Job_type_Full Time", "Job_type_Intern", etc.

- Title, Place, and Department have been target encoded using the "Overall_rating" as the target. This means the values in these columns represent the average overall rating for each category in

13

the original columns. For instance, the "Title_encoded" column contains average overall ratings for each job title.
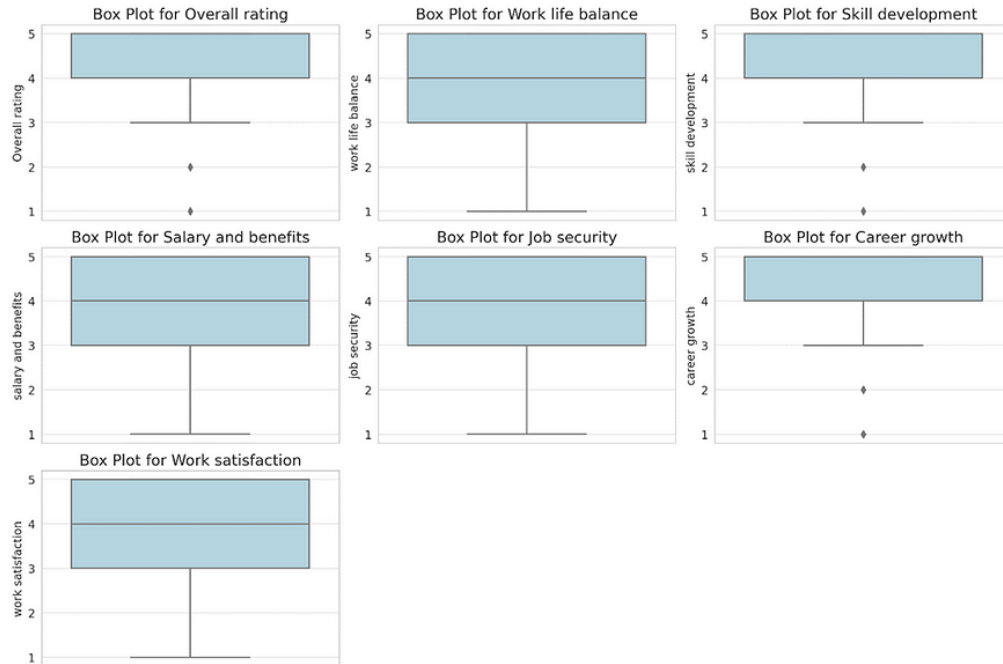
**Handling outliers:**

Outliers are extreme values that can skew our data and potentially bias our model's estimates. We'll focus on the numerical rating columns, as these are most likely to contain outliers.

The steps involved are:

1.  Visualize potential outliers using box plots.

2.  Decide on a strategy to handle the outliers (e.g., removing, capping, or transforming).

Let's start by visualizing the potential outliers using box plots for the rating columns.



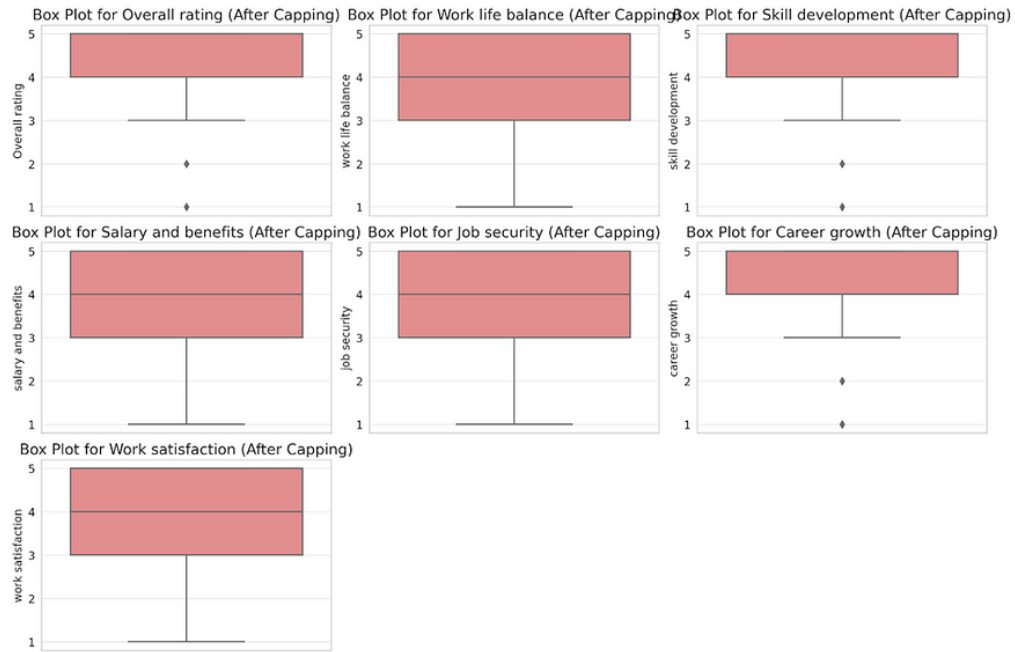The box plots help visualize the spread and potential outliers for each rating column:

- Overall_rating, work_life_balance, skill_development, salary_and_benefits, and job_security have some lower outliers.

- Career_growth and work_satisfaction seem to have both lower and upper outliers, though the number is small.

Given the nature of the data (ratings), we should be cautious about removing outliers because they might represent genuine feedback. Instead, we can consider the following strategies:

1.  Capping: Replace outlier values with a certain threshold (e.g., 1st or 99th percentile values).

2.  Transformation: Apply transformations like log or square root to reduce the impact of outliers.

For this dataset, capping might be more appropriate since applying transformations like log to rating data may not be intuitive.

Let's cap the outliers at the 1st and 99th percentiles for each rating column. This approach ensures that extreme values are limited to a certain threshold, reducing their potential influence without entirely removing the observations.



The box plots, post-capping, show that the outliers have been effectively limited to the 1st and 99th percentiles for each rating column. The color difference highlights the post-capping state.

By capping the outliers rather than removing them, we retain the original observations and reduce their extreme influence.

Before we transition to modeling, it's crucial to ensure our data is properly prepared. Feature scaling is an essential step, especially for algorithms that rely on distances or gradients.

For this dataset, we'll apply Min-Max scaling to the numerical columns to ensure all features are on a similar scale.

```
from sklearn.preprocessing import MinMaxScaler

# Columns to scale
columns_to_scale = rating_columns + ['Likes_Sentiment', 'Dislikes_Sentiment',
'Review_Month', 'Review_Year']

# Applying Min-Max scaling
scaler = MinMaxScaler()
data_encoded[columns_to_scale] = scaler.fit_transform(data_encoded[columns_to_scale])

# Checking the first few rows post scaling
data_encoded[columns_to_scale].head()
```

| Overall_rating | work_life_balance | skill_development | salary_and_benefits | job_security | career_growth | work_satisfaction | Likes_Sentiment | Dislikes_Sentiment | Review_Month | Review_Year |
|---|---|---|---|---|---|---|---|---|---|---|
| 1.00 | 0.5 | 0.50 | 0.75 | 1.00 | 0.50 | 0.50 | 0.5000 | 0.500000 | 0.727273 | 1.0 |
| 1.00 | 1.0 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 0.9775 | 0.727500 | 0.727273 | 1.0 |
| 0.75 | 1.0 | 0.25 | 0.25 | 0.75 | 0.50 | 0.00 | 0.5000 | 0.500000 | 0.727273 | 1.0 |
| 1.00 | 1.0 | 1.00 | 1.00 | 1.00 | 0.75 | 1.00 | 0.7500 | 0.500000 | 0.727273 | 1.0 |
| 0.75 | 0.0 | 0.00 | 0.50 | 0.75 | 0.00 | 0.25 | 0.8500 | 0.291667 | 0.727273 | 1.0 |

The features have been successfully scaled using Min-Max scaling. Now, all the numerical columns are within the range [0, 1], making them ready for most machine learning models.

## SEMMA: Model

In the "Model" phase, we'll focus on building and evaluating models to capture patterns in the dataset. Given the nature of our data (ratings and feedback), a regression model might be suitable to predict the "Overall_rating" based on other features. Additionally, clustering can provide insights into different segments of reviews.

Our steps will include:

1.      Splitting the data into training and testing sets.

2.      Building a regression model.

3.      Evaluating model performance.

4.      (Optional) Clustering to identify segments.

Let's start by splitting the data into training and testing sets to ensure we can evaluate our model's performance on unseen data.

The data has been split into training and testing sets:

- Training data: 8,601 observations

- Testing data: 2,151 observations

With our data ready, let's move on to building a regression model. For simplicity and efficiency, we'll start with a Linear Regression model. This model will attempt to predict the "Overall_rating" based on the provided features.

```
lr_model.fit(X_train, y_train)

# Predicting on the test set
y_pred = lr_model.predict(X_test)

# Evaluating the model performance
mse = mean_squared_error(y_test, y_pred)
rmse = mean_squared_error(y_test, y_pred, squared=False)
r2 = r2_score(y_test, y_pred)
```

The Linear Regression model has been trained and evaluated:

1. Mean Squared Error (MSE): 0.04640.0464

2. Root Mean Squared Error (RMSE): 0.21550.2155

3. R-squared ($R^2$): 0.40710.4071

The $R^2$ value indicates that the model explains approximately 40.71% of the variability in the "Overall_rating". While not exceptionally high, this is a decent starting point given the complexity and inherent subjectivity of review data.

To further enhance our model's performance, let's explore an ensemble technique: Random Forest Regression.

Random Forest is a versatile algorithm that can work well for both classification and regression tasks. By building multiple decision trees and averaging their predictions, it can capture complex relationships in the data.

Here are the steps we'll follow:

1. Train a Random Forest Regressor.

2. Evaluate its performance on the testing set.

3. Compare its performance with the Linear Regression model.

```
from sklearn.ensemble import RandomForestRegressor

# Initializing and training the Random Forest Regressor
rf_model = RandomForestRegressor(n_estimators=100, random_state=42)
rf_model.fit(X_train, y_train)
```

```
# Predicting on the test set
y_pred_rf = rf_model.predict(X_test)

# Evaluating the model performance
mse_rf = mean_squared_error(y_test, y_pred_rf)
rmse_rf = mean_squared_error(y_test, y_pred_rf, squared=False)
r2_rf = r2_score(y_test, y_pred_rf)
```

The Random Forest Regressor's performance is as follows:

1. Mean Squared Error (MSE): 0.02430.0243

2. Root Mean Squared Error (RMSE): 0.15590.1559

3. R-squared ($R^2$): 0.68970.6897

Comparatively, the Random Forest Regressor has a higher $R^2$ value of approximately 68.97%, indicating it explains nearly 69% of the variability in the "Overall_rating". This is a significant improvement over the Linear Regression model.

## SEMMA: Assess

In the "Assess" phase, the key goal is to evaluate the performance, robustness, and reliability of our models. We've already touched on this earlier when we evaluated the Linear Regression and planned to evaluate the Random Forest Regressor models. However, assessment goes beyond just model metrics.

Here's what we typically consider in this phase:

1. Model Performance Metrics: As we've done, use metrics like $R^2$, MSE, RMSE, etc., for regression problems to understand how well our model is performing.

2. Model Robustness: Evaluate how the model performs under different conditions or with different subsets of the data. This can involve techniques like cross-validation.

3. Business Impact: Consider the real-world implications of the model's predictions. For instance, what's the impact of a false positive versus a false negative in our context? This often requires collaboration with domain experts.

4. Residual Analysis: For regression models, analyzing residuals (differences between observed and predicted values) can give insights into areas where the model might be underperforming. Patterns in residuals can indicate that the model isn't capturing all the underlying dynamics in the data.

5. Model Interpretability: Understand which features are driving the model's predictions. This is especially important for complex

models or when communicating results to stakeholders. Tools like SHAP (SHapley Additive exPlanations) or feature importance from tree-based models can help here.

## Deployment and Monitoring

After assessing and finalizing our models, the next steps in a real-world scenario would typically be:

1. Deployment:

- API Integration: Models can be deployed as APIs, allowing applications or other systems to make real-time predictions.

- Batch Predictions: For non-real-time predictions, the model can be run on a batch of data periodically.

- Tools: Tools like Flask, FastAPI, or platforms like AWS Sage-Maker, Azure ML, or Google AI Platform can be used for deployment.

2. Monitoring:

- Model Drift: Over time, the model's performance can degrade as the underlying data distributions change. Monitoring tools can help detect this drift.

- Feedback Loop: Implementing a feedback loop allows users or systems to provide feedback on the model's predictions, enabling continuous improvement.

- Retraining: Based on the feedback and monitoring, the model can be retrained periodically with new data.

3. Reporting:

- Dashboards: Tools like Tableau, Power BI, or even libraries like Dash for Python can be used to create dashboards to visualize the model's performance and other KPIs.

- Alerts: Set up automated alerts for drastic performance drops or other anomalies.

4. Iterative Development:

- Feature Engineering: Based on model insights and feedback, new features can be engineered or existing ones refined.

- Model Selection: Periodically, it's good practice to revisit the model selection step to see if other algorithms or newer techniques provide better performance.

5. Stakeholder Communication:

- Keeping stakeholders updated on the model's performance, potential improvements, and any changes is crucial. This can involve periodic meetings, reports, or dashboards.

6. Documentation:

- Comprehensive documentation ensures that the data science process, decisions, and models are transparent, replicable, and maintainable by others in the future.

# Conclusion

We embarked on this journey with the objective of analyzing the Tata Motors Employee Reviews dataset using the SEMMA methodology. The dataset provided valuable insights extracted from employee reviews, capturing sentiments, ratings, and other relevant details about the company's work environment.

**Sample:**

After understanding the dataset's nature, we found that it contained multiple features, including ratings for various aspects of the job, textual reviews, job titles, and more. The primary challenge was to handle missing values and derive insights from a mix of numerical and textual data.

**Explore:**

Through Exploratory Data Analysis (EDA), we discovered:

- The distribution of various ratings, identifying that most ratings were above average.

- Trends in the number of reviews over time, which indicated periodic peaks.

- Sentiment scores derived from textual reviews, highlighting the general sentiment trend in "Likes" and "Dislikes" sections.

**Modify:**

Data preprocessing included:

- Handling missing values, primarily through median imputation.

- Feature engineering, including extracting sentiment scores from reviews and encoding categorical variables.

- Capping outliers for rating columns to minimize their impact.

- Scaling features to prepare the dataset for modeling.

**Model:**

We explored regression models to predict the "Overall_rating" based on other features. Our journey covered:

- A Linear Regression model, which explained around 40.71% of the variability in the "Overall_rating".

- A Random Forest Regressor, which showed a significant improvement, explaining nearly 69% of the variability.

- Theoretically discussed clustering, intending to group similar reviews and gain insights into different review segments.

**Assess:**

We assessed the model's performance and robustness. Continuous monitoring, feedback loops, and iterative development were emphasized as essential for real-world applications. We also touched upon the importance of model interpretability.

**Real-World Implications & Considerations:**

The insights derived from the data analysis can guide interventions and strategic decisions at Tata Motors. Ethical considerations, data privacy, scalability, and stakeholder buy-in were highlighted as crucial aspects when transitioning from analysis to implementation.

## Final Thoughts:

The Tata Motors Employee Reviews dataset offers a glimpse into employees' perceptions and experiences. Through systematic data analysis using SEMMA, organizations can identify areas of improvement, strategize interventions, and continually adapt to enhance employee satisfaction. As with any data-driven approach, it's essential to balance quantitative insights with qualitative understanding, ensuring that actions taken resonate with the workforce's genuine needs and concerns.

## References

SAS Institute Inc. (2008). SEMMA Data Mining Methodology. Cary, NC: SAS Institute Inc.

AmbitionBox (2023). Tata Motors Employee Reviews. Retrieved from AmbitionBox Website.

Breiman, L. (2001). Random forests. Machine learning, 45(1), 5-32.

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). An Introduction to Statistical Learning. Springer.

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Vanderplas, J. (2011). Scikit-learn: Machine learning in Python. Journal

of machine learning research, 12(85), 2825-2830.

Pang, B., & Lee, L. (2008). Opinion mining and sentiment analysis. Foundations and Trends® in Information Retrieval, 2(1–2), 1-135.

Bird, S., Klein, E., & Loper, E. (2009). Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit. O'Reilly Media, Inc.

McKinsey & Company. (2017). Driving value from advanced analytics in HR. McKinsey & Company Insights.