# Statistical Analysis of Job and Wage Trends with Clustering Analysis

**Introduction**

The project aims to analyze and uncover trends in wages, job growth, educational demands, and sectoral employment patterns across various industries. This analysis seeks to identify critical insights into pre- and post-COVID salary changes, popular job roles, and emerging trends in employment and education.

Given the vast scope and complexity of the data, clustering analysis was chosen as a core methodology to segment occupations and industries based on wage ranges, educational qualifications, and regional demands. Additionally, machine learning techniques were employed to project future trends and create a career recommendation system.

**Objective of the Analysis**

The primary goals of this project include:
- Identifying popular occupations and the evolving demands across sectors.
- Highlighting salary trends before and after the COVID-19 pandemic.
- Understanding educational requirements for specific job roles and their alignment with wages.
- Clustering occupations into meaningful groups to simplify analysis and draw actionable insights.
- Developing a recommendation system for job seekers to align education levels with high-demand occupations and competitive salaries.
- Forecasting job and wage trends for the next five years using machine learning techniques.

**Why Clustering Analysis Was Chosen?**

Clustering analysis is a powerful unsupervised machine learning technique that helps identify natural groupings in the data without pre-defined labels. It was chosen for this project because:

- *Complex Dataset*: The dataset contains multiple attributes such as wages, education levels, regional data, and job trends, which are difficult to analyze through simple descriptive statistics.
- *Pattern Discovery*: Clustering allows us to identify hidden patterns and segment occupations into meaningful clusters.
- *Dimensionality Reduction*: Clustering simplifies the analysis of large, multi-dimensional datasets by reducing them into smaller, interpretable groups.
- *Actionable Insights*: Clustering helps in grouping jobs with similar wage patterns or educational demands, which is useful for trend analysis and decision-making.

**Why Other Clustering Techniques Were Not Used?**

Several clustering techniques were considered before finalizing the optimal approach for this analysis. Here is a brief evaluation:

- *Hierarchical Clustering*: While this technique provides a visual hierarchy of clusters, it is computationally expensive for large datasets and is less suitable for identifying optimal clusters in high-dimensional data.
- *DBSCAN (Density-Based Clustering)*: DBSCAN performs well for data with varying densities but struggles when the dataset does not exhibit clear density variations, as in our case.
- *Gaussian Mixture Models (GMM)*: GMM assumes a Gaussian distribution of data points, which did not align well with the structure of our dataset.

After careful consideration, K-Means Clustering was chosen due to its simplicity, computational efficiency, and effectiveness in handling numerical, high-dimensional data.

**Methodology to Obtain Optimal Clustering Results**

To ensure the best clustering outcomes, the following methodologies and steps were implemented:

*Step 1: Data Preprocessing*

- Data Cleaning: Missing values were imputed, and redundant data was removed to improve data quality.
- Feature Selection: Key features such as wage ranges, education levels, job trends, and regional demands were selected based on domain relevance.
- Data Normalization: Features were normalized using Min-Max scaling to bring all variables to a similar scale, ensuring fair distance computations.

*Step 2: Choosing the Optimal Number of Clusters*

To identify the optimal number of clusters, the following methods were used:

*Elbow Method*: The Within-Cluster Sum of Squares (WCSS) was plotted against the number of clusters, and the point of inflection (elbow) was chosen as the optimal number.

*Silhouette Analysis*: The silhouette score, which measures the cohesion and separation of clusters, was used to validate the results.

Both methods consistently suggested an optimal range for the number of clusters (e.g., k = 3 or 4).

*Step 3: K-Means Clustering Implementation*

- K-Means Clustering was applied using Python's `sklearn` library.
- The algorithm iteratively assigned data points to clusters based on distance from centroids and recalculated cluster centroids until convergence.
- Final clusters represented groups of occupations with similar wage trends, educational demands, or regional patterns.

*Step 4: Cluster Interpretation*
**Cluster Analysis**: Each cluster was analyzed based on its centroid values to identify the characteristics of jobs and sectors within it.
**Visualization**: Clusters were visualized using techniques like PCA (Principal Component Analysis) to project high-dimensional data into a 2D space for interpretability.

**Results of Clustering Analysis**
The clustering analysis produced the following key results:

1. *Cluster 1: High-Wage Occupations*
- Jobs requiring advanced degrees (Master's/Ph.D.)
- Prevalent in technology, finance, and healthcare sectors
- Post-COVID wage increases observed in these roles.

2. *Cluster 2: Medium-Wage Occupations*
- Jobs requiring bachelor's degrees or diplomas
- Common in manufacturing, education, and public administration
- Moderate wage growth post-COVID with regional variations.

3. *Cluster 3: Low-Wage Occupations*
- Jobs requiring minimal education (high school or below)
- Prevalent in retail, hospitality, and logistics sectors
- Significant wage stagnation post-COVID.

**Key Insights**
1. *Pre- and Post-COVID Trends*: High-wage roles in technology and healthcare saw significant growth post-COVID, while low-wage sectors experienced stagnation.
2. *Educational Demands*: Advanced degrees align with high-paying roles, whereas entry-level roles offer lower wages.
3. *Regional Variations*: Urban regions demonstrate stronger demand for high-wage roles compared to rural areas.

**Conclusion**

Clustering analysis proved to be an effective technique for segmenting occupations and identifying wage trends across sectors. By implementing K-Means Clustering and validating results through the Elbow Method and Silhouette Analysis, the analysis uncovered actionable insights into job growth, wage patterns, and educational demands.

This analysis not only helps job seekers align their education with emerging opportunities but also aids policymakers and stakeholders in addressing wage disparities and regional demands.

**Future Scope**

- *Integration with Time-Series Analysis*: Forecast wage and job trends for specific clusters over the next decade.
- *Career Recommendation System*: Enhance the existing recommendation model using content-based filtering to provide personalized job and education suggestions.
- *Exploration of Advanced Clustering Techniques*: Implement algorithms like DBSCAN or GMM for further refinement of results.
- *Interactive Dashboards*: Create visual dashboards for stakeholders to interact with the analysis and explore cluster-specific trends.

This comprehensive clustering analysis delivers valuable insights into the evolving landscape of occupations, education demands, and wage trends, empowering decision-makers with data-driven recommendations.

# Career Recommendation System with Content-Based Filtering

**Objective of the Analysis**
The primary purpose of this project is to help individuals make informed career choices based on their education level and desired salary. By leveraging data-driven insights, the system recommends jobs that align with the user's education qualifications and salary expectations.

**Why Content-Based Filtering Was Chosen?**

1. *Personalized Recommendations*: Content-based filtering recommends careers based on the user's provided features, such as education and salary, making it ideal for tailored suggestions.

2. *Interpretability*: Content-based methods use user input and combine it with job-related attributes. This provides transparency into why a specific career was recommended.

3. *Sufficient Feature Data*: The dataset includes valuable attributes like job title, education level, work experience, and salary, which can be effectively combined to create meaningful "features."

4. *No Historical User Data Needed*: Unlike collaborative filtering, content-based filtering doesn't rely on historical user behavior or feedback, making it suitable for first-time users or systems with no prior user data.

**Methodology Used -**

**1. Data Preprocessing:**
- Salary data was binned into ranges such as `salary_40k`, `salary_50k`, etc., for better representation in the recommendation model.
- Missing values in "Work Experience" were filled with "No prior experience."
- Combined features were created using `Typical Education`, `Work Experience`, and `Salary` as input for TF-IDF.

**2. Feature Engineering:**
The Combined_Features column was constructed to represent each job description in a simplified text format.
For example: `"Bachelor's degree No prior experience salary_50k"`

**3. TF-IDF Vectorization:**
- **Why TF-IDF?** TF-IDF (Term Frequency-Inverse Document Frequency) measures the importance of words in the combined text features while reducing the impact of frequently occurring generic words.
- The job descriptions and user inputs were converted into TF-IDF vectors for numerical representation.

**4. Cosine Similarity:**

   - Used to compute the similarity between user-provided input (education and desired salary) and the available job features.

  - Jobs with the highest similarity scores were recommended to the user.

**Alternative Recommendation Techniques and How They Can Be Incorporated -**

1. **Collaborative Filtering:**

*Approach*:

Instead of relying on job attributes, collaborative filtering recommends jobs based on user behavior, preferences, or feedback.

Example: If users with similar qualifications have applied for a specific job role, this role could be recommended to new users with matching profiles.

*Integration*:
- Collect user behavior data such as job clicks, preferences, or feedback ratings over time.
- Implement matrix factorization techniques (e.g., Singular Value Decomposition - SVD).

*Advantages*:
- Works well when sufficient user interaction data is available.
- Can recommend jobs even if feature information (like salary or work experience) is incomplete.

2. **Hybrid Recommendation System:**

*Approach*:

Combine content-based filtering & collaborative filtering to leverage both job attributes & user behavior.

*Integration*:
- Use content-based filtering for cold-start users who lack historical interaction data.
- Gradually incorporate collaborative filtering as users provide feedback or interact with the system.

*Advantages*:
- Mitigates limitations of individual methods (e.g., cold-start problem, lack of content diversity).
- Provides a more robust and personalized recommendation.

3. **Clustering-Based Recommendation (K-Means Clustering)**:

*Approach*: Group jobs into clusters based on similar attributes (e.g., salary, work experience, education).

*Integration*:
- Use K-Means or DBSCAN clustering to group jobs into categories.
- Recommend jobs from the closest cluster based on user preferences.

*Advantages*:
- Efficient for handling large datasets.
- Provides job category insights that can help users explore similar careers.

4.   **Deep Learning Techniques**:
*Approach*: Use deep neural networks to learn job and user representations (embeddings).
*Example*: Use techniques like Neural Collaborative Filtering (NCF) or Autoencoders to recommend jobs.
*Integration*:
- Build embeddings for jobs based on features like education, work experience, and salary.
- Train the model with user feedback (ratings or clicks) to improve recommendations over time.
*Advantages*:
- Handles complex relationships between job attributes.
- Highly scalable for large datasets.

5.   **Knowledge-Based Recommendation**:
*Approach*: Use predefined rules and domain knowledge to recommend jobs.
*Example*:
If a user selects a "Master's degree" with a desired salary of `$80,000`, recommend managerial roles in specific industries.
*Integration*:
Define industry-specific rules based on qualifications, salary ranges, and work experience.
*Advantages*:
- Useful when there's domain-specific information.
- Offers straightforward and explainable recommendations.

**Why TF-IDF and Cosine Similarity Over Other Techniques?**

1. *Simplicity and Efficiency*:
TF-IDF combined with Cosine Similarity is computationally efficient, making it ideal for systems with limited user interaction data.
2. *Transparency*:
It is easy to explain the recommendation process to users, as it directly uses their input to find the most relevant jobs.
3. *Cold-Start Problem*:
Unlike collaborative filtering, content-based filtering does not require historical user interaction data, making it perfect for systems in the early stages.
4. *Scalability*:
The model can scale with large datasets, as vectorization and similarity computation can handle high-dimensional data efficiently.

**Future Enhancements -**

1. *Incorporating User Feedback*:
- Allow users to provide feedback on the relevance of job recommendations.
- Use feedback to retrain models (e.g., collaborative filtering) for improved results.
2. *Integration of NLP Techniques*:
Use advanced NLP models like BERT or Sentence Transformers for richer textual understanding and better similarity calculations.
3. *Geographical Constraints*:
Add location-based filtering to recommend jobs specific to the user's region.
4. *Dynamic Salary Insights*:
Incorporate real-time salary trends to ensure recommendations align with market conditions.
5. *Visualization*:
Include graphs or charts showing job demand trends based on user qualifications and salary expectations.

**Conclusion -**
The current implementation of the **Career Recommendation System** effectively uses content-based filtering with TF-IDF and cosine similarity to recommend jobs tailored to users' education levels and salary preferences. By incorporating alternative methods like collaborative filtering, clustering, or deep learning, the system can evolve into a more robust and intelligent career recommendation engine.

# Occupational Projection: Techniques for Job Trends Prediction

**Random Forest for Occupational Projection**
Random Forest is a widely-used ensemble machine learning technique for both regression and classification tasks. It is particularly well-suited for predicting future job trends due to its robustness, accuracy, and ability to handle complex relationships between features.

*Key Advantages of Random Forest*:
• *Handles Complex Relationships*: Random Forest can model intricate, non-linear relationships between features such as industry type, education level, and job growth rates.
• *Feature Importance*: It automatically computes feature importance, allowing identification of the most critical factors influencing job trends, such as education level, region, or industry type.
• *Robust to Overfitting*: By aggregating predictions from multiple decision trees trained on random subsets of data, Random Forest reduces overfitting and improves generalizability.
• *Handles Missing Data*: Missing values are effectively handled by averaging predictions or selecting the best splits during training, which is particularly useful for incomplete datasets in job-related surveys or labor reports.
• *Flexibility*: Random Forest can process both numerical (e.g., wages, employment growth) and categorical data (e.g., job titles, education levels) without requiring extensive preprocessing.

**How Random Forest Works:**
- A Random Forest consists of multiple decision trees trained independently on random subsets of the training data.
- For regression tasks, such as predicting the employment growth rate or median annual wage for occupations, the model averages the predictions of all trees.
- For classification tasks, such as determining whether a job will grow or decline, the model uses a majority vote from all trees.
- Random Forest introduces randomness by (1) bootstrapping data (sampling with replacement) and (2) randomly selecting features at each split, which ensures diversity among trees and improves performance.

**Insights Gained**:
- *Key Drivers of Job Growth*: Feature importance scores will highlight the key predictors of occupational trends (e.g., regional demand, education levels, post-pandemic recovery).
- *Industry-Specific Trends*: Random Forest can reveal which industries are projected to experience the most growth or decline.
- *Education-Level Impacts*: It can predict which education levels align with future high-demand occupations, providing guidance for career planning.

**2. Gradient Boosting Methods (XGBoost, LightGBM, CatBoost)**
Gradient Boosting techniques are highly effective for predictive modeling tasks, providing accurate projections for job trends.

**Why Use Gradient Boosting?**
- *Superior Accuracy*: Gradient Boosting optimizes predictions by sequentially improving weak learners (decision trees).
- *Handles Imbalanced Data*: It performs well even when certain jobs or sectors dominate the dataset.
- *Feature Importance*: Similar to Random Forest, Gradient Boosting computes feature importance to pinpoint key factors driving occupational trends.
- *Customization*: Modern libraries like XGBoost, LightGBM, and CatBoost allow tuning of hyper-parameters to improve performance and reduce training time.

**Real-World Application:**
- Use XGBoost for high-dimensional datasets where numerous variables influence job projections.
- Use LightGBM for faster computations on large datasets.
- CatBoost is ideal for datasets with categorical variables (e.g., education levels, job types) without extensive encoding.

**Time Series Forecasting (ARIMA, SARIMA, and Prophet)**
Time series forecasting techniques are crucial for predicting job growth or decline over time.

**Key Techniques:**
1.  *ARIMA (Auto-Regressive Integrated Moving Average)*:
- Suitable for linear trends and seasonal data.
- Can model employment data and forecast trends based on historical time-series data.

2.  *SARIMA (Seasonal ARIMA)*:
- Extends ARIMA by incorporating seasonality, making it ideal for capturing periodic fluctuations in employment trends.

3.  *Prophet (by Facebook)*:
- User-friendly for handling irregular data and missing values.
- Can handle holidays, seasonal effects, and trend shifts (e.g., post-COVID impacts on certain industries).

*Application Example*:
Forecast the employment growth of industries over the next 5 years based on historical job market trends and seasonal patterns.

4.  *Neural Networks (LSTM and GRU for Sequential Data)*
Long Short-Term Memory (LSTM) and Gated Recurrent Units (GRU) are types of Recurrent Neural Networks (RNNs) capable of learning sequential dependencies in time series data.

**Why Use LSTMs and GRUs?**
- Captures Long-Term Dependencies: LSTM and GRU can model complex, non-linear trends in time series data without losing long-term dependencies.
- Non-Stationary Data: Ideal for occupational projections where data may exhibit changing trends due to economic or social factors.
- Dynamic Predictions: Unlike traditional methods, these neural networks adapt to changes in job trends over time.

Application:
Use LSTM models to predict employment trends, wage fluctuations, and job demand across regions.

5. *Clustering Techniques (K-Means, DBSCAN, Hierarchical Clustering)*

Clustering techniques help group occupations with similar characteristics, revealing hidden patterns and trends.

Methods:

*K-Means Clustering*:
- Group occupations based on features like median wages, education level, and employment growth.
- Identify clusters of high-demand jobs or declining industries.

*DBSCAN (Density-Based Clustering)*:
- Detects outlier occupations with unusual growth or decline patterns.

*Hierarchical Clustering*:
- Creates a hierarchical structure of job categories to explore relationships between job types.

**Insights**: Cluster analysis can identify "hot spots" of job growth or decline across sectors, guiding job seekers and policymakers.

6. *Regression Models (Linear and Polynomial Regression)*

While simpler than ensemble and deep learning models, regression techniques are still valuable for projecting job trends.

- Linear Regression: Estimates the relationship between features (e.g., education level, industry demand) and employment projections.
- Polynomial Regression: Captures non-linear trends by extending linear regression to higher-degree polynomial features.

Example Application: Predict how changes in education levels affect median annual wages for specific occupations.

7. *Natural Language Processing (NLP) for Job Descriptions*

NLP techniques, such as TF-IDF and BERT, can analyze job descriptions to identify emerging skills and trends in the labor market.

Techniques:
- TF-IDF (Term Frequency-Inverse Document Frequency): Extracts important keywords from job descriptions to determine trending skills.
- BERT (Bidirectional Encoder Representations from Transformers): Provides advanced language understanding to analyze job postings.

Use Case: Identify which skills and qualifications are becoming more critical for future job opportunities.

**Conclusion**

By incorporating a combination of advanced techniques such as Random Forest, Gradient Boosting, Time Series Forecasting, Neural Networks, Clustering, and Regression Models, we can build a robust occupational projection system. Each technique offers unique strengths, allowing us to forecast employment trends, identify growth opportunities, and guide stakeholders such as policymakers, job seekers, and educators.

Key factors like education levels, regional demand, industry-specific trends, and historical data can be leveraged to generate actionable insights for the evolving job market. Combining ensemble methods, time series models, and NLP ensures the system is versatile, accurate, and aligned with real-world data complexities.

**Code for Occupational Projection: Key Considerations**

When building a model for occupational projection, the code typically integrates multiple stages, from data preprocessing to model deployment.
Here's a structured overview of what the implementation involves:
1. *Data Collection and Preprocessing*
The first step is gathering relevant data, such as:
• Historical Employment Data (e.g., government labor reports, industry trends).
• Education Levels and Demographics (from surveys or databases).
• Economic Indicators (GDP, inflation, post-pandemic recovery).
• Regional Factors (state-wise job trends, urban vs. rural demand).

Once the data is collected, preprocessing techniques are applied:
- *Data Cleaning*: Handle missing values, outliers, and inconsistencies.
- *Feature Engineering*: Create meaningful features, such as growth rates, industry codes, or region categories.
- *Encoding*: Convert categorical variables (e.g., education levels, job titles) into machine-readable formats using one-hot encoding or label encoding.
- *Normalization/Scaling*: Standardize numerical features (like wages or growth percentages) for better model performance.

2. Model Selection and Implementation:
The choice of model depends on the type of projection (classification or regression).
• Random Forest: Suitable for capturing non-linear relationships and identifying important predictors of job trends.
• Gradient Boosting (e.g., XGBoost, LightGBM): Enhances accuracy by sequentially learning from previous errors.
• Time Series Forecasting (ARIMA, Prophet): Predicts trends in employment growth or job demand over time.
• LSTM/GRU: Neural networks model complex sequential patterns in time-series data.

*Implementation Workflow*:
- Split data into training, validation, and testing sets.
- Train models using frameworks like scikit-learn, XGBoost, or PyTorch.
- Optimize hyper parameters using techniques such as GridSearchCV or Random Search.
- Evaluate performance using metrics like RMSE(for regression) or accuracy/F1-score (for classification).

3. *Insights and Feature Importance*

Once the model is trained, the next step is interpreting the results:
- *Feature Importance*: For Random Forest or Gradient Boosting, visualize the most influential features driving projections, like education levels or industries.
- *Model Interpretation*: Use tools like SHAP (SHapley Additive exPlanations) or LIME to understand how features contribute to predictions.

These insights help policymakers, educators, and job seekers make informed decisions about skill development and industry opportunities.

4. *Deployment and Visualization*

Finally, the model can be deployed using:
- *APIs*: Build endpoints using tools like Flask or FastAPI for real-time predictions.
- *Visualization Tools*: Use libraries like Matplotlib, Seaborn, or Plotly to create dashboards showcasing job trends, projections, and key growth areas.
- *Scalability*: Deploy on cloud platforms like AWS, Azure, or Google Cloud for scalable occupational projections across regions.

By combining robust modeling with clear visualization, the code for occupational projection ensures accuracy and usability for diverse stakeholders.