

Project Report: AI-Powered CV Analysis for Career Guidance

1. Problem Statement

In today's competitive job market, university students and unemployed graduates often face significant challenges in understanding job market expectations and optimizing their resumes. Two primary pain points were identified:

1. **Salary Expectation Alignment:** Job seekers frequently lack clear insight into appropriate salary ranges for specific roles, leading to uncertainty when applying for positions. This project aims to provide a data-driven salary estimation to assist individuals in assessing a job's salary potential *before* committing to an application.
2. **Resume Fortification & Benchmarking:** Students and graduates struggle to create resumes that meet professional standards. The project addresses this by developing a system to evaluate a resume's professionalism and completeness against a benchmark of professional resumes, providing a quantitative score (between 0 and 1). A score above 0.6 is deemed indicative of a strong, professional resume, offering actionable insights for improvement.

The overarching goal was to streamline aspects of the recruitment process and empower job seekers with intelligent tools for career guidance and resume optimization.

2. Data

The project utilized two distinct datasets to achieve its objectives:

- **Resume Classification Dataset:** This dataset was sourced from a GitHub project, containing information on over 13,000 individuals, primarily in the form of raw resume text. This dataset served as the foundation for training the CV classification model, allowing it to learn patterns and features associated with various professional categories.
- **Salary Estimation Dataset:** While not explicitly detailed as a separate input for model training during our collaborative code development, the intention was to leverage salary data, likely sourced from platforms like Glassdoor (via Kaggle), to inform the salary estimation component. This component would conceptually map resume features (e.g., experience, skills, education) to predicted salary ranges.

Data Preprocessing: A robust text preprocessing pipeline was implemented to prepare the raw CV text for machine learning:

- **Text Cleaning:** All CV text was converted to lowercase, punctuation was removed, numerical digits were stripped, and multiple whitespace characters were normalized to single spaces. This ensured consistency and reduced noise.

- **Text Vectorization:** A `tf.keras.layers.TextVectorization` layer was employed to convert the cleaned text into numerical sequences suitable for model input. The key parameters for this layer were set as:
 - `max_tokens = 15000`: Limiting the vocabulary to the top 15,000 most frequent words.
 - `output_sequence_length = 512`: Ensuring all vectorized CV sequences had a uniform length, padding or truncating as necessary.
- **Label Encoding:** Resume categories were transformed from categorical text labels into numerical integers using `sklearn.preprocessing.LabelEncoder`.

3. Model Architecture

The core of this project involved a deep learning model, primarily focusing on the CV classification task, with a conceptual integration for salary estimation.

CV Classification Model: A Bidirectional Long Short-Term Memory (Bi-LSTM) neural network architecture was chosen for its effectiveness in processing sequential data like text, capturing long-range dependencies.

- **Embedding Layer:**
 - `input_dim: MAX_FEATURES (15,000 + 1 for OOV tokens)`
 - `output_dim: 256` (This dense vector representation helps the model learn semantic relationships between words.)
- **Bidirectional LSTM Layer:**
 - `units: 128` (Processes sequences in both forward and backward directions, enhancing context understanding.)
- **Dense Layers:**
 - A sequence of `Dense` (fully connected) layers with `relu` activation functions to learn complex patterns from the LSTM outputs. The exact number and size of these layers were tuned for optimal performance.
- **Output Layer:**
 - A final `Dense` layer with `softmax` activation, with the number of units equal to the total number of unique resume categories. This layer outputs a probability distribution over all possible categories.

Salary Estimation Component (Conceptual Integration): While a distinct regression model architecture was not explicitly detailed during our shared code development, the project's scope included salary estimation. This capability would typically be integrated either as:

- A **multi-task learning head** attached to the common feature extraction layers of the CV classification model, using a separate regression output layer (e.g., a `Dense` layer with no activation for a continuous output).
- A **separate regression model** trained on extracted features from the CV text (such as years of experience, specific skills, education level, implied seniority from job titles), mapping these to a numerical salary value. Evaluation for this would involve regression-specific metrics.

4. Evaluation

The primary evaluation focused on the performance of the CV classification model.

CV Classification Performance:

- **Loss Function:** `SparseCategoricalCrossentropy` was used, which is suitable for multi-class classification problems where labels are integers.
- **Metrics:** The model's performance was evaluated using `accuracy`. While exact final values depend on the specific training run, the goal was to achieve high training and validation accuracy, indicating the model's ability to generalize well to unseen resumes. The iterative refinement of the model architecture and training parameters aimed to consistently push these metrics higher.

Salary Estimation Performance: For the salary estimation component, typical evaluation metrics for regression tasks would include:

- **Mean Absolute Error (MAE):** Represents the average absolute difference between predicted and actual salary values.
- **Root Mean Squared Error (RMSE):** Provides a measure of the magnitude of errors, with larger errors given disproportionately more weight. Actual MAE/RMSE values would depend on the implementation and training of the salary estimation module.

5. Project Lifecycle Stages

The project followed a standard machine learning development lifecycle, characterized by iterative refinement and problem-solving:

1. **Data Acquisition & Understanding:** Sourcing and initial exploration of the resume and salary datasets.
2. **Data Preprocessing & Cleaning:** Implementing robust functions for text cleaning, tokenization, vectorization, and label encoding to transform raw data into a machine-readable format.
3. **Model Design & Development:** Conceptualizing and implementing the deep learning architecture, starting with an initial design and progressively enhancing it (e.g., adding Bidirectional LSTM layers, adjusting neuron counts).
4. **Model Training:** Training the classification model on the prepared dataset, monitoring loss and accuracy metrics over epochs.
5. **Model Evaluation:** Assessing the model's performance on a validation set to ensure generalization and identify areas for improvement.
6. **Model Persistence:** Implementing robust saving mechanisms for the trained model (using the `.keras` format) and preprocessing assets (like the `TextVectorization` vocabulary and `LabelEncoder`) to ensure reproducibility and ease of deployment.
7. **Prediction Implementation:** Developing the code to load the saved model and preprocessing tools, allowing for real-time inference on new, unseen CVs.

8. **Iterative Debugging & Refinement:** Continuously identifying and resolving technical errors and performance issues encountered at various stages.

6. Challenges & Solutions

Throughout the project development, several technical challenges were encountered and systematically addressed, demonstrating robust problem-solving:

- **InvalidArgumentError (Vocabulary Mismatch):** This error occurred during prediction, indicating that the `TextVectorization` layer was producing indices outside the range expected by the `Embedding` layer of the loaded model.
 - **Solution:** This was resolved by ensuring that the `MAX_FEATURES` parameter used during the `TextVectorization` layer's training/adaptation (and subsequent saving) was perfectly synchronized with the `MAX_TOKENS` (or `input_dim`) parameter of the `Embedding` layer within the Keras model during both training and loading. Retraining and saving the model with the updated `MAX_FEATURES` was crucial.
- **AttributeError: 'TFSMLayer' object has no attribute 'predict' (Model Loading Method):** This issue arose when attempting to use the `.predict()` method on a model loaded via `TFSMLayer`.
 - **Solution:** The model saving and loading strategy was revised to exclusively use the Keras `.keras` format (`model.save('model.keras')` and `tf.keras.models.load_model('model.keras')`). This ensured that the loaded object was a standard Keras model with the `.predict()` method readily available.
- **ValueError: File not found (Pathing Issues):** Persistent errors related to files not being found during model or asset loading.
 - **Solution:** Meticulous verification of file paths, ensuring correct Google Drive mounting, and using Colab's built-in "Copy path" feature to eliminate typos. The consistency of `SAVE_DIR` and `TF_MODEL_LOAD_PATH` was paramount.
- **Initial Model Performance & Specific Predictions (e.g., 'Automobile' for Software Engineer):** Early iterations of the classification model sometimes made unexpected predictions, indicating a need for architectural improvement.
 - **Solution:** The model architecture was enhanced, notably by incorporating `Bidirectional LSTM` layers. This significantly improved the model's ability to understand context within the text, leading to more accurate classifications across diverse categories. Increasing `MAX_FEATURES` also allowed the model to learn from a richer vocabulary.

These challenges, though common in machine learning projects, were effectively overcome through systematic debugging and applying appropriate TensorFlow/Keras best practices.

7. Conclusion & Future Work

This project successfully developed and demonstrated an AI-powered system for CV analysis, capable of classifying resumes into professional categories and conceptually integrating salary estimation. By leveraging deep learning techniques and robust data preprocessing, the system provides valuable insights for both job seekers (through resume evaluation and salary benchmarking) and recruiters (by streamlining initial resume screening).

Key Achievements:

- Built a functional deep learning model for multi-class CV classification with good accuracy.
- Implemented a comprehensive text preprocessing pipeline.
- Overcame significant technical challenges related to model saving, loading, and text vectorization.
- Provided a framework for integrating salary estimation based on resume content.

Future Work & Potential Improvements:

- **Enhance Salary Estimation:** Fully develop and integrate a dedicated regression model for salary estimation, potentially incorporating more granular features (e.g., explicit years of experience parsing, specific technology skill matching).
- **Expand Dataset:** Incorporate a larger and more diverse dataset for both classification and salary estimation to further improve model robustness and generalization.
- **Fine-tuning & Advanced Architectures:** Explore more sophisticated deep learning architectures (e.g., Transformer-based models like BERT) or advanced fine-tuning techniques for even higher accuracy.
- **Explainability (XAI):** Implement methods to explain *why* a resume was classified into a certain category or received a particular score, providing more actionable feedback for users.
- **User Interface & Deployment:** Develop a user-friendly web application or API to allow seamless interaction with the model for real-time CV analysis.
- **Multi-task Learning Refinement:** Explore multi-task learning architectures where the classification and salary estimation tasks share common deep features, potentially leading to more efficient and accurate learning for both.