# A Two-Step Resume Information Extraction Algorithm

### 1.1 Motivation

This paper addresses the challenge of extracting information from resumes with diverse formats by proposing a two-step information extraction algorithm. The key contribution lies in the introduction of a feature called "Writing Style" that captures the inherent structure within each line of the resume. This method aims to improve the accuracy of information extraction and reduce reliance on manual labelling compared to existing approaches.

### 1.2 Methodology

The proposed algorithm consists of two steps:

Text Block Classification: This step categorises each line of text in the resume into different sections (e.g., Education, Work Experience) without relying on a predefined template. It utilises a multi-class classifier to identify potential block titles (e.g., "Education") based on word frequency and synonym clustering.

Resume Facts Identification: This step focuses on identifying specific details within each section. It leverages the "Writing Style" feature, which combines information about word classification, punctuation, and phrase location within the line. The paper mentions using a cosine similarity measure and clustering for entity recognition.

The algorithm is evaluated on a real-world dataset of 15,000 Chinese resumes

### 1.3 Conclusion

The paper concludes that the proposed method achieves promising results in terms of precision and recall compared to existing methods based on HMM or CRF. This approach offers advantages in reducing manual labelling efforts and potentially handling diverse resume formats.

## Limitations

### 2.1 Limited Implementation Details

i.The paper lacks in-depth explanations of certain aspects.

ii.Data cleaning processes for handling raw text extracted from resumes are not detailed.

iii.Specific algorithms used for text block classification and named entity recognition are not mentioned.

iv.The evaluation methodology for text block classification relies on results from other research papers, which limits the transparency of the overall evaluation.

### 2.2 Lack of Baseline Comparison

While the paper justifies not comparing with rule-based methods due to their limitations, it would be beneficial to include a comparison with a simpler information retrieval baseline to strengthen the evaluation and showcase the relative improvement achieved by the proposed method.

## Synthesis

### 3.1 Potential Applications

The accurate extraction of information from resumes can be valuable for various applications, including:

Building resume repositories for recruitment purposes.

Automating the process of candidate screening and selection.

Populating applicant tracking systems with extracted information.

### 3.2 Future Directions

The paper proposes interesting future work directions:

Exploring techniques like multi-label prediction and coreference resolution to improve named entity recognition.

Investigating extensions of the "Writing Style" feature to handle more complex layouts or multilingual resumes.
Utilising the extracted information to discover social relations among people (although the paper doesn't elaborate on specific methods or applications).
These future directions hold promise for further enhancing the accuracy and applicability of the proposed information extraction method.

**Name :Ikti Safat Anjum Soumya**
**ID: 21201816**