# Algorithmic Segmentation of Job Ads Using Textual Analysis

Benjamin Murauer[1(✉)], Michael Tschuggnall[1], Günther Specht[1], and Julia Brandl[2]

[1] Department of Computer Science, University of Innsbruck, Innsbruck, Austria
`b.murauer@posteo.de`, `michael.tschuggnall@ches.pro`,
`guenther.specht@uibk.ac.at`
[2] Department of Organization and Learning, University of Innsbruck, Innsbruck, Austria
`julia.brandl@uibk.ac.at`

**Abstract.** As job ads are getting more prevalent online, an automated analysis is becoming increasingly important, especially in the field of human resource management. In this paper, we propose an approach to automatically segment job ads by predefined categories like the description of a job or the offering company, which is needed to categorize and quantify different aspects of job ads. Using a manually annotated data set, textual features are extracted for each segment type in a first step and utilized to train state-of-the-art machine learning classification methods. Subsequently, these models are used by iterative algorithms to detect the individual segments. Using several optimization techniques like detecting typical segment start phrases, comprehensive evaluations show promising results.

**Keywords:** Text segmentation · Job ads analysis · Structured text classification

## 1 Introduction

With the rise of the interconnected world wide web, the amount of data digitally transferred through various channels is increasing steadily. Especially businesses have adjusted and optimized their workflows by utilizing new possibilities of data exchange in many areas like allowing location-separated global teams or performing targeted social media marketing campaigns. Consequently, also the publication of job ads has been expanded or even completely shifted from print media to online services. With massive amounts of digitally available job ads, a crucial task in the field of human resource management (HRM) is to systematically analyze them to be able to answer questions like "What does the market want?", "Which companies search for which person types?" or "In which domains are social skills more important?" [4,21].

Typical study examples include attempts to detect trending key requirements demanded by employers in a specific domain, or the extraction of what is offered

by companies in return (e.g., [1,3]). Most of those studies are done by manual inspection of advertisements, which has the drawback that only a small subset of potential sources can be examined due to limitations of human resources. To conduct quantitative research, an automated textual analysis is needed. A fundamental prerequisite for making meaningful statements is thereby an automated identification of well-known segments like the description of the job tasks or requirements within ads. For example, it makes a crucial difference if the phrase "work experience" is stated within the job requirements or appears in a sentence like "You will be integrated in a team with long-term work experience", as different conclusions can be drawn depending on the position of the phrase.

Given that a job ad depicts the employer's offer and the future employee's duties, the job ad content can be divided into four categories [6], whereby not necessarily all of them have to be present: the (i) **company description** introduces the offering company, who they are, what they do etc., the (ii) **job description** describes the job itself with all its components, in the (iii) **requirements** it is stated which education, technical/social skills etc. are demanded from the applicant, and finally the (iv) **offer** segment describes what the company offers, including salary and other benefits. Examples of those four categories are shown in Table 1.[1]

This paper presents an approach to automatically detect those categories, if present, by applying textual analysis based on state-of-the-art machine learning techniques. It can thereby be seen as a hybrid approach that incorporates both machine learning techniques as well as specifically developed algorithms operating on top. Respectively, the two major steps are: (i) training a model in order to be able to classify each segment, and (ii) calculating the final segmentation utilizing this model. In the first step, a data set with manually assigned ground truth is used to train a model which is able to identify each segment type with high accuracy when it is given correctly separated segments as input, i.e., one of the four categories and not overlapping structures. Using this model, algorithms are designed to enable an automatic classification even if the borders are not previously known, and thus segmenting and clustering the job ad by the defined categories.

Summarizing the main contribution of this paper, we present an algorithmic text segmentation approach that is specifically tailored for job ads. Utilizing job-ad-specific knowledge, experiments reveal that algorithmic segmentation represents a valid and comparable alternative to Conditional Random Fields (CRF) [17], a state-of-the-art technique used in text segmentation and tagging.

The rest of this paper is organized as follows: Sect. 2 briefly outlines related work, and Sect. 3 explains the generation of the model in detail, which is needed as a prerequisite for the actual job ad segmentation presented in Sect. 4. The resulting basic algorithm is evaluated in Sect. 5, which is improved and reevaluated in Sect. 6. Finally, a conclusion and possible future work is discussed in Sect. 7.

---

[1] Note that while the example is provided in English for better understandability, the data used throughout this paper is in German.

**Table 1.** Examples for the different segment types (anonymized).

| Category | Example 1 | Example 2 |
|---|---|---|
| Company description | In 2017, *CompanyX* not only launched a completely new product, it created and has led ever since a whole new product category. Nowadays *CompanyX* employs more than 12000 people in over 165 countries, selling over 5 billion products a year. The World of *CompanyX* provides the forum for you to use your talent and passion, to develop yourself and make an impact | About us: we look for innovation everywhere. For 130 years, we have been at the forefront of innovation, but finding solutions to the world's biggest problems has never been more important than right now. Join us today and become an essential part of the solution! |
| Job description | The Controller role is responsible for acting as business partner to the Sr. VP or VP of the defined area and the respective management team. Project lead and participation in the respective functional area is also part of the role. In addition, steering support processes within his/her defined area and contributing to core financial processes is expected from the position holder as well | Your responsibilities<br>– Work on the software verification (component test, system test, building verification environment)<br>– Develop automated test cases and work closely with the development team to ensure testability<br>– Take part in the analysis of problems that occur in live operations<br>– Travel in order to directly support our clients in the verification phase of the project |
| Requirements | Because our business is dynamic and advances in science and technology require new methods of production we are looking for individuals who can do the following: (i) collaborate with team members in the identification and implementation of continuous improvement initiatives and action plans. (ii) Support activities in the areas of cost containment, efficiency, productivity, energy conservation, waste minimization, operational excellence and lean practices | Your profile<br>– Experience with Scrum or other iterative methodologies and experience with C# Development<br>– Minimum of 4 years of experience designing and developing systems in a .Net environment<br>– Experience with JavaScript, not only jQuery<br>– Prior experience with a source control system (SVN, GIT, TFS, . . . )<br>– Good verbal and written communication skills in English and German (at least level B2) |
| Offer | The minimum annual salary for the position is 50,000$ gross, whereby the effective salary is based on your qualification and experience | We offer:<br>– an interesting job within the packaging & paper industry<br>– to be part of a successful multicultural company<br>– an empowering environment<br>– new office building in the city center<br>– several attractive social benefits |

## 2   Related Work

Job advertisements have become more prevalent online, and they have shown to improve the chances of being employed [16]. Therefore, research on job ads is becoming more important in many fields, with classic examples ranging from optimizing skill sets for specific positions [1] to optimizing strategies for HR departments [11]. The information to be gained from the text differs greatly. For example, for some applications detailed features that are required by potential applicants are extracted [3], while others try to extract information to fill specific templates [7]. Also, complex frameworks have been constructed to automatically use job ads and distribute them over different channels [15], where detailed information is extracted by rule-based models, including regular expressions and other manually tailored methods. In order to improve the extraction procedures, one

goal is to detect specific topics within the text, whereby there exists no general rule of how many parts a job should be divided into. For example, [18] extract 13 different categories from the text of French job ads, including, e.g., contact information or mobility requirements, whereas [34] divide their documents into 3 more general sections (education, job title and sector). Both approaches rely on common textual features like word, character or POS-tag frequencies.

To cope with the general text and/or topic segmentation problem, a wide range of different methods is in use, often based on the research by Hearst [14], in which the lexical cohesion of terms is analyzed. Generally, approaches utilizing the lexical cohesion of segments follow either a *local* or *global* strategy. Local methods thereby process documents using fixed-sized fragments, aiming to locally detect significant changes in cohesion (e.g., [12,14]). On the other hand, global methods try to maximize the lexical cohesion for each segment from a global perspective (e.g. [8,32]). Due to the fact that both strategies have their pros and cons, i.e., leading to different accuracies in different situations, e.g., Simon et al. propose a graph-based hybrid approach [29]. With respect to this terminology, the presented approach in this paper can also be seen as a text segmentation problem that falls into the hybrid category, as it locally traverses text blocks using sliding windows while considering the global view. Nevertheless, in contrast to typical approaches in this field, the ad segmentation does not aim to draw borders based on topic or genre changes, but by predefined job ad categories. Moreover, segments must be attributed to those categories, i.e., only finding borders is insufficient in this case.

In terms of style markers used, approaches utilize different features and methods, e.g., Bayesian models [10], Hidden Markov Models [5], vocabulary analysis in various forms like word stem repetitions [24] or word frequency models [25]. While there exist some recent papers (e.g., [20,26]) that include a comparison between some of the segmentation approaches on the same data sets, in general it is difficult to compare performances due to the heterogeneous problem and data types being approached. To counter this problem, the PAN workshop series[2] recently proposed tasks to create stylistic clusters within documents, revealing that accurate stylistic segmentation is still a difficult problem [30,31].

With respect to automatically labelling tokens from a text document (e.g., POS tagging or Named Entity Recognition), Conditional Random Fields (CRF) have proven to be effective [19,28]. Moreover, they have successfully been used for other sequential labeling problems, such as detecting common fields like authors, title, year or publisher from the headers and citations of research papers [22].

## 3 Preliminary Experiments: Training a Textual Model

To tackle the problem of automatically detecting the four predefined segment types (i) *company description*, (ii) *job description*, (iii) *requirements* and (iv) *offer* within an unseen job ad, at first a textual model is needed that can differentiate between those segments. I.e., arbitrarily given one of those segments, the model should automatically be able to predict the correct one based on textual

---

characteristics. In this work, two variants, i.e., models based on dictionaries and machine-learned models have been built and tested. Both variants are presented in the following and are created on the basis of a manually tagged data set consisting of approximately 1,200 manually tagged Austrian job ads, written in German language[3].

### 3.1   Dictionary-Based Models

Due to the fact that job ads are formulated very heterogeneously, e.g., containing full sentences describing the job in detail, or consisting basically only of a job title accompanied by bullet lists, the main unit taken into account are words (including stemming and n-grams). As a first attempt, dictionaries for each category are created using the tagged data set. The prediction is then made by comparing the given segment with the dictionaries and computing similarities. Concretely, the following methods have been tested:

(1.) *Simple Stems:* In this variant the dictionary of a segment consists of all occurring stems and their normalized weights. Given an unseen segment $S$, the similarity to a dictionary $D$ is then calculated by summing up the weights in $D$ for every occurring stem in $S$. By using this computation, important stems that frequently appear in a specific segment are also of higher importance with respect to the final similarity score.

(2.) *N-Grams:* This approach is very similar to the previous one and only replaces stems by character n-grams (using $n = \{2, 3, 4, 5, 6\}$). Thus, a dictionary contains all occurring n-grams including their normalized weights, and the similarity is computed by summing up the weights of all matching n-grams.

(3.) *Lucene:* Finally, the query model of the Java library *Lucene*[4] has been utilized, which uses the tf-idf metric and vector space models. Thereby a document is created for each segment type by concatenating all samples, which is subsequently given to Lucene to create an appropriate index. To be specific, four different indices are created, each consisting of only one document containing all samples of the respective job ad category. Finally, the similarity of an unseen segment $S$ to a dictionary $D$ is then computed as follows: For each word in $S$, a query is formed that consists of this word only, and the similarity score for the dictionary document is retrieved from Lucene. The overall score is then computed by summing up the delivered scores of all queries (words).

### 3.2   Machine-Learned Models

As an alternative to the dictionary-based prediction, common machine learning techniques have been applied to build a model which can predict the segment type. The list of provided features is as follows: (i) single word stems

---

*(1s)*, stemmed word 2-grams *(2s)* and stemmed word 3-grams *(3s)*, each of them by removing or keeping stop words (*nonstop/stop*), (ii) all possible combinations of the previous features (e.g., *1s-3s-nostop* for using single stems and stemmed word 3-grams, by eliminating stop words), and (iii) character n-grams using $n = \{2, 3, 4, 5, 6\}$.

To determine the best working algorithm for this approach, several commonly used methods have been tested. Using the WEKA toolkit as a general framework [13], the following classifiers have been utilized: Naive Bayes a Bayes Network using the K2 classifier, Large Linear Classification using LibLinear, a support vector machine using LibSVM with nu-SVC classification, a k-Nearest-Neighbours classifier (kNN) using $k = 1$, PART, and a pruned C4.5 decision tree (J48). Thereby all classifiers have been used with their respective standard settings.

## 3.3 Model Evaluation

Both the dictionary-based as well as the machine-learned models have been evaluated using a manually tagged data set, which has been created from a subset of 1,200 samples of all Austrian job ads in 2015. All considered ads are formulated in German language, and manually segmented into the four segment types by an expert group of three persons with high experience in the human resource management field. The content of the samples thereby varies from mainly just bullet lists to fully formulated, grammatically correct sentences. Moreover, segment types may also be spread over multiple positions in the text, e.g., stating a job requirement at the beginning as well as at the end of an ad.

The data set has been divided into two halves, whereby one half has been used as the training/test corpus for the evaluation of the models, and the other half for the evaluation of the segmentation algorithm (see Sect. 5). From the 600 job ads serving for the model evaluation, 300 have been used to train the models, i.e., to build the dictionaries or to train the classifiers, respectively, and the remainder has been used for testing. Finally, the size of the dictionaries ($d_s$) has been parameterized (except for Lucene), choosing the dictionary to contain only the 50, 75, 100, 150 or 200 most frequent stems/n-grams, or to contain all of them. In case of the machine-learned models, the number of features has also been parameterized by the same criteria.

**Table 2.** Best accuracies of dictionary-based and machine-learned models in percent, given already segmented categories.

| Method | $d_s$ | cmp | Job | Offer | Req. | Avg. |
|---|---|---|---|---|---|---|
| Lucene | all | 85.7 | 84.8 | 89.2 | 93.1 | **88.2** |
| 5-grams | 75 | 73.2 | 76.6 | 81.4 | 88.3 | **79.9** |
| 6-grams | 75 | 72.3 | 76.6 | 81.8 | 84.0 | **78.7** |
| stems | 200 | 75.3 | 72.7 | 77.9 | 87.0 | **78.2** |
| 4-grams | 75 | 69.3 | 75.3 | 74.9 | 81.4 | **75.2** |
| 3-grams | 100 | 61.0 | 69.3 | 66.7 | 75.8 | **68.2** |

(a) Dictionary-Based Models

| Classifier | $d_s$ | Features | cmp | Job | Offer | Req | Avg |
|---|---|---|---|---|---|---|---|
| LibSVM | 100 | 1s-2s-stop | 94.8 | 94.7 | 97.0 | 97.9 | **96.1** |
| BayesNet | 150 | 4-grams | 92.5 | 92.7 | 97.2 | 96.1 | **94.6** |
| Naive Bayes | 75 | 1s-3s-stop | 93.0 | 92.1 | 95.5 | 96.4 | **94.2** |
| kNN | 150 | 5-grams | 88.2 | 92.1 | 93.2 | 96.8 | **92.6** |
| LibLinear | 150 | 4-grams | 88.0 | 93.1 | 96.3 | 90.8 | **92.0** |
| PART | 100 | 1s-3s-stop | 81.0 | 84.8 | 90.4 | 89.4 | **86.4** |
| J48 | 100 | 1s-3s-stop | 79.1 | 82.2 | 91.2 | 89.6 | **85.5** |

(b) Machine-Learned Models

The best results for both dictionary-based and machine-learned models are shown in Table 2, where it can be seen that – while Lucene performs best for the former – most machine learning algorithms outperform it substantially. By utilizing frequencies of single stems and stemmed 2-grams, LibSVM could achieve an accuracy of 96%. In general, good accuracies are achieved for all four segment types by restricting the number of features, with the *offer* and *requirement* type having slightly better values.
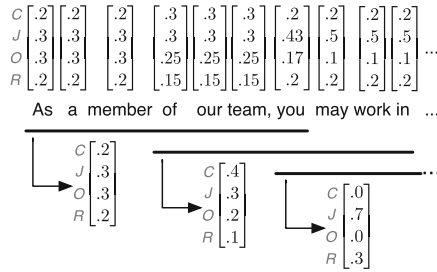
## 4    Ad Segmentation

According to the nature of most job ads, which are given as unstructured full texts, the previously discussed classification cannot be used as is, because the borders of the different segments are not known. That is, a classifier can only predict the segment type with high accuracy if the complete and correctly separated segment is given as input. To tackle this limitation, an algorithm has been developed that is based on the constrained classifier, but allows to estimate segment borders. In the first step, a model-based probability is calculated for each text position to belong to each segment type. Subsequently, these values serve as input for the final ad segmentation, basically by locating probability peaks and applying thresholds to widen a segment.

### 4.1    Applying the Model

Given a full text job ad, at first a probability for each segment type and text position (token) is estimated. This is done by iteratively traversing the text using sliding windows of length $w_l$ and step $w_s$, using the model to predict a vector $[p_C, p_J, p_O, p_R]$ for each window. The values in this vector correspond to the probability for the window to be of the type *company*, *job*, *offer* and *requirement*, respectively. For example, Fig. 1 shows three windows of length seven, where for each of the windows a corresponding vector is computed using the model. The window step $w_s$ in the example is three tokens.

After processing all sliding windows, the probability vector for each token is computed by calculating the average of all predicted window vectors where the respective token appears. In the example, the first three words have been predicted by only one window, and thus the probability vectors for those words correspond to the vector of the corresponding window. On the other side, the tokens *'of'*, *'our'*, *'team'* have been predicted by two windows, and thus the token vectors represent the average of those. Finally, *'you'* has been predicted by three windows and is consequently calculated as the average of all three involved windows.

The overall procedure can be formalized as is stated in Algorithm 1. Given the consecutive list of tokens $T$ of the job ad, it computes the average probability vector $P_i$ for each token $T_i$ by applying the model function $prob(t)$ on the respective tokens $t$ of each window. According to the results of Sect. 3.3, machine-learned models are used for calculating the respective probabilities. Thereby, as

**Fig. 1.** Example of computing probability vectors for the categories company description (C), job description (J), offer (O) and requirements (R).

not all classifiers deliver an estimation for each class (segment type) but only the predicted class (e.g., LibSVM would only produce $[0, 0, 1, 0]$, whereas BayesNet would produce a more informative vector like $[.1, .1, .5, .3]$), the model function $prob(t)$ utilizes BayesNet instead of the slightly more accurate LibSVM classifier.

A visual example of the probability calculation for a sample job ad is depicted in Fig. 2, which shows the probabilities for each token and segment type. The text document spans on the x-axis from left to right, and the ground truth for each type is included in colored rectangles. It can be seen that each segment type is matched by the probability distribution, having peaks in the correct corresponding segment.[5]

## 4.2   Calculating the Final Segmentation

On the basis of the probability vectors for each token and segment type, the aim of the last step is to compute the final segmentation. As can be seen visually, probabilities like in Fig. 2 already match their corresponding segment type, i.e., having their peaks within the correct frame. To determine the whole fragments algorithmically, the basic idea is to use thresholds that widen a peak to the left and the right, respectively, until the probability falls below the thresholds. According to Algorithm 2, the procedure is as follows: (1.) Out of the unattributed segment types, choose the one with the highest peak. (2.) Starting from the peak, go to the left and attribute every token to the chosen type, until the threshold is reached or the token is already attributed to another type. (3.) Starting from the peak, do the same to the right.

During the execution, it may happen that peaks for a segment type appear within an already attributed segment. In this case, this segment type has no attribution, which also reflects the possibility that a job ad may not contain a

---

[5] What can be observed in addition is that not all tokens correspond to a segment type, i.e., the text between company and job description and at the end of the ad is general and thus not attributed to a specific type.

certain segment type. In a similar way, segment types have no final attribution if an unattributed peak is found, but left and right positions are already below the given thresholds.

**Algorithm 1 .** Probability Vector Calculation

**input:**

| | |
|---|---|
| $T$ | list of tokens |
| $w_l, w_s$ | window length, window step |
| $prob(t)$ | computes the probability vector for the tokens $t$ |

**output:**

| | |
|---|---|
| $P$ | list of average probability vectors |

1: **for** $i$ from 1 to length$(T)$ **do** $v_i \leftarrow []$
2: **end for**

3: **for** $i$ from 1 step $w_s$ to length$(T)$ **do**
4:     $t = [T_i, ..., T_{i+w_l}]$
5:     $[p_C, p_J, p_O, p_R] \leftarrow prob(t)$
6:     **for** $j$ from $i$ to $i + w_l$ **do**
7:         append $[p_C, p_J, p_O, p_R]$ to $v_j$
8:     **end for**
9: **end for**

10: **for** $i$ from 1 to length$(T)$ **do**
11:     $a_C \leftarrow$ average of all $p_C$ in $v_j$
12:     $a_J \leftarrow$ average of all $p_J$ in $v_j$
13:     $a_O \leftarrow$ average of all $p_O$ in $v_j$
14:     $a_R \leftarrow$ average of all $p_R$ in $v_j$
15:     $P_i \leftarrow [a_C, a_J, a_O, a_R]$
16: **end for**

**Algorithm 2.** Final Segmentation

**input:**

| | |
|---|---|
| $P_t$ | token probabilities for type $t$ |
| $\delta_l, \delta_r$ | left and right thresholds |
| $maxP(T)$ | returns the segment type with the highest peak |

**output:**

| | |
|---|---|
| $S$ | segment type attribution per token |

1: **for** $i$ from 1 to length$(P)$ **do** $S_i \leftarrow$ *unknown*
2: **end for**
3: $types \leftarrow \{C, J, O, R\}$

4: **while** $|types| > 0$ **do**
5:     $t \leftarrow maxP(types)$
6:     $p \leftarrow$ index of peak within $P_t$,     $i \leftarrow p$
7:     **while** $P_{ti} \geq \delta_l \wedge S_i = unknown$ **do**
8:         $S_i \leftarrow t$,     $i \leftarrow i - 1$
9:     **end while**
10:     $i \leftarrow p + 1$
11:     **while** $P_{ti} \geq \delta_r \wedge S_i = unknown$ **do**
12:         $S_i \leftarrow t$,     $i \leftarrow i + 1$
13:     **end while**
14:     $types \leftarrow types - \{t\}$
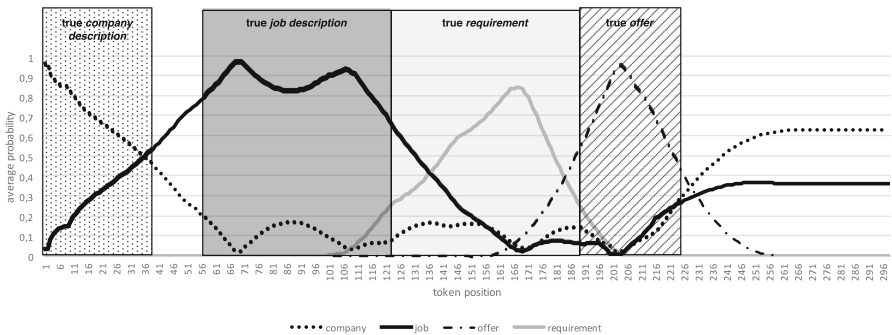15: **end while**



**Fig. 2.** Example of segment type probabilities for a sample job ad.

## 5   Evaluation

The presented approach has been evaluated using the remainder of the data set described in Sect. 3.3, i.e., using 600 manually tagged Austrian job ads. According to the input of the algorithms, the following parameter ranges have been evaluated: window length in tokens $w_l = \{15, 25, 50, 75\}$, window step in tokens $w_s = \{1, 2, ..., 10\}$, left and right thresholds $\delta_l = \delta_r = \{0.3, 0.4, ..., 0.8\}$. Typically, studies on job ads focus on analyzing the appearances of individual tokens within different segment types[6] rather than requiring exact border positions (e.g., [2,21]). Therefore, recall, precision and $F_1$-score have been used as the main metrics, which are calculated based on the comparison between the correct tokens of a segment and the predicted tokens[7]. With this method, the exact locations of segment borders are only implicitly evaluated, as only tokens - regardless of their text position - are taken into account. To also get evaluation results with respect to border positions, the commonly used text segmentation metric *WindowDiff* [23] as well as its slightly altered variant *WinPR* [27] have been utilized in addition. WindowDiff yields a normalized value between 0 and 1, where 0 indicates a perfect result, whereas WinPR computes standard recall and precision values for computed segmentations. Note that WindowDiff as well as WinPR only measure the accuracy of border positions and do not reflect the correct attribution of segments.

To put the results in perspective, two baselines have been computed: (i) the *RND-BASELINE* simply divides the job ad into four equally long segments and assigns the respective categories in a random order, and (ii) *CJRO-BASELINE* also divides the text equally into four segments, but always uses the most prominent category order for attribution, which is: *company description – job description – requirements – offer*.

Table 3 shows the best results per window length. Using a length of 25, i.e., 25 tokens per window, a token-based F-score of 67% could be gained, significantly outperforming the baselines. The segmentation-based F-score of WinPR is at 53% only, which indicates that the accuracy of finding exact border positions is substantially worse than determining relevant tokens within segments.

**Table 3.** Evaluation results of the basic segmentation algorithm.

| $w_l$ | $w_s$ | $\delta_l$ | $\delta_r$ | segments | | | | overall | | | winDiff | winP | winR | winF |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | $F_{(C)}$ | $F_{(J)}$ | $F_{(O)}$ | $F_{(R)}$ | recall | prec | F | | | | |
| 25 | 1 | 0.65 | 0.35 | 65.7 | 73.8 | 54.0 | 71.5 | 65.2 | 70.2 | **67.6** | 35.0 | 52.8 | 53.9 | 53.3 |
| 15 | 1 | 0.6 | 0.3 | 62.7 | 75.8 | 46.6 | 67.5 | 62.2 | 70.2 | **66.0** | 35.0 | 50.4 | 50.7 | 50.5 |
| 50 | 1 | 0.65 | 0.35 | 67.2 | 61.6 | 62.7 | 70.1 | 70.1 | 61.8 | **65.7** | 37.1 | 50.8 | 47.0 | 48.8 |
| 75 | 1 | 0.65 | 0.35 | 66.0 | 56.3 | 64.3 | 65.6 | 70.6 | 57.1 | **63.1** | 38.4 | 48.7 | 43.4 | 45.9 |
| CJRO-BASELINE | | | | 62.1 | 56.1 | 30.5 | 41.0 | 61.7 | 41.6 | 49.1 | 44.2 | 39.4 | 19.7 | 26.1 |
| RND-BASELINE | | | | 25.3 | 21.9 | 21.0 | 22.4 | 29.6 | 20.0 | 23.4 | 44.2 | 39.4 | 19.7 | 26.1 |

---

[6] E.g., if 'abroad' appears within the offer and/or requirement segment.

[7] Thereby, if a segment type does not exist in the job ad but at least one token has been predicted for it, the recall/precision is set to 0, and also vice versa.

## 6   Improvements

Based on the previous results from Sect. 5, several techniques have been applied in order to further improve the accuracy of the approach.

1. An intrinsic characteristic of job ads is that bullet lists are used frequently, e.g., to express tasks or required skills. Moreover, it seems reasonable that segment types do not switch within bullet lists[8]. Therefore the algorithm is extended to be able to also detect bullet lists using regular expressions. If a bullet list is detected during the spanning of segments to the left and right from the probability peak, the algorithm steps over the whole bullet list and includes it in the current segment. By doing so, it is ensured that bullet lists are never split, regardless if the probability is below the left or right threshold or not.

2. When manually inspecting job ads, it can be observed that individual segment types often start with similar phrases. For example, the description of a company may start with *'About us:'*, *'Who we are:'* or *'[COMPANY_NAME] is a ...'*. In order to detect such typical start sequences, additional dictionaries have been created using the training data, which store all stemmed token sequences of segment beginnings. If a newline character is found within the first three tokens, only this line is used, otherwise the first three tokens are used. Having computed a dictionary for every segment type, the segmentation algorithm from Sect. 4.2 is then altered to utilize segment beginnings. This means, that before searching for the probability peak within a segment type, it is first assessed whether a typical start sequence could be found. If yes, the matching position is used as the beginning of the segment, from where the algorithm spans now only to the right (using $\delta_r$) and not to the left. To decide if a start position is found, common regular expressions on stems are used which match exactly, i.e., with no tolerance. On the other hand, if also the company name is involved, the extracted name of the offering company given by the data set is utilized. Because the extracted name often does not correspond exactly to the name written in the job ad, fuzzy string matching is applied[9].

Both previously mentioned modifications have been evaluated and reveal that both can enhance the performance when applied individually[10]. By finally combining both improvement methods, an accuracy of 77% could be gained, as can be seen in Table 4. To compare the performance to existing segmentation approaches that do not only draw borders but are also capable of assigning classes to segments, a conditional random field (CRF) with standard parameter configuration has been trained. It was then utilized to label every token with one of the four categories, plus an artificial *"other"* label for tokens belonging to no category. It can be seen that the presented algorithmic approach outperforms the

---

[8] Although there are rare cases where this is the case.
[9] Using the Jaro-Winkler distance [9,33] with a threshold of 0.85.
[10] Due to space constraints, the individual results are omited.

**Table 4.** Evaluation results including bullet list and start sequence detection.

| $w_l$ | $w_s$ | $\delta_l$ | $\delta_r$ | Segments | | | | Overall | | | winDiff | winP | winR | winF |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | $F_{(C)}$ | $F_{(J)}$ | $F_{(O)}$ | $F_{(R)}$ | recall | prec | F | | | | |
| 50 | 1 | 0.5 | 0.35 | 69.9 | 79.0 | 72.8 | 86.6 | 79.6 | 75.3 | **77.4** | 28.5 | 63.9 | 57.4 | 60.5 |
| 75 | 1 | 0.55 | 0.45 | 69.4 | 78.6 | 72.9 | 88.1 | 79.4 | 75.6 | **77.4** | 29.2 | 62.6 | 57.2 | 59.8 |
| 25 | 1 | 0.4 | 0.35 | 69.3 | 73.5 | 68.3 | 73.6 | 70.2 | 74.1 | **72.1** | 29.9 | 61.9 | 59.4 | 60.6 |
| CRF-BASELINE | | | | 64.1 | 71.0 | 73.1 | 72.1 | 71.9 | 71.6 | **71.5** | – | – | – | – |
| 15 | 1 | 0.4 | 0.3 | 67.8 | 76.5 | 60.3 | 69.7 | 67.3 | 74.0 | **70.5** | 30.2 | 61.0 | 56.3 | 58.5 |

CRF-baseline[11]. Note that while in principle it would be possible to incorporate at least some of the algorithmic improvements also to the CRF-baseline, this task is not trivial as the CRF provides no probabilities but only labels tokens. Therefore, a different and adapted approach would have to be developed which is not addressed in this paper and left for future work.

In a final set of experiments using the hybrid approach, it has additionally been evaluated whether introducing a minimum segment length $l_{min}$ increases the overall accuracy. By defining $l_{min} = 28$ tokens, the overall F-score could be improved slightly to 78%. On the other side, an experiment was conducted that allowed a segment to be distributed over several positions within the document[12]. This was done by spanning a segment not only from one probability peak, but from two or three peaks if an individual predefined threshold was exceeded for the latter. Optimizing this threshold it has been found that the overall accuracy could not be improved, but instead is reduced by 10%. In a final experiment, it has been evaluated whether allowing different left and right thresholds for different segment types influences the segmentation accuracy, but the results indicated that this is not the case.

## 7   Conclusion and Future Work

In this paper, we presented an approach to automatically segment job ads into four predefined segment types. Using machine-learning based on extracted textual features, iterative algorithms operate to detect segments. Comprehensive evaluations including several optimizations using domain specific characteristics reveal that the segment types can be identified with an accuracy of 78% in average, outperforming an out-of-the-box Conditional Random Field implementation as a baseline. The proposed method may also be applied to other languages than German and evaluated, how the accuracy changes with different data sets. Moreover, it should be evaluated whether training a model which includes not only the predefined segment types, but also a specific general type (*'other'*) can enhance the accuracy. Finally, the classifier that holds the segment type model should be tuned for an application in real-life business scenarios.

---

[11] Note that the segmentation metrics winDiff and winPR are not applicable for the CRF-baseline, as tokens are only labeled and no actual segmentation is performed.

[12] Which is the case for the *offer* segment, for example, if the salary is stated within the *job description*, while the other "offer" is formulated at a different position.

# References

1. Aken, A., Litecky, C., Ahmad, A., Nelson, J.: Mining for computing jobs. IEEE Softw. **27**(1), 78–85 (2010)
2. Anonymous Authors. Anonymized Title. Working Paper (2017)
3. Aureli, E., Iezzi, D.F.: Recruitment via web and information technology: a model for ranking the competences in job market. In: Proceedings of JADT 2006, pp. 79–88 (2006)
4. Jansen, B.J., et al.: Using the web to look for work: implications for online job seeking and recruiting. Internet Res. **15**(1), 49–66 (2005)
5. Blei, D.M., Moreno, P.J.: Topic segmentation with an aspect hidden Markov model. In: Proceedings of the 24th ACM SIGIR Conference on Research and Development in Information Retrieval, New York, NY, USA, pp. 343–348 (2001)
6. Budd, J.W., Bhave, D.: The employment relationship. In: The SAGE Handbook of Human Resource Management, pp. 51–70 (2010)
7. Califf, M.E., Mooney, R.J.: Bottom-up relational learning of pattern matching rules for information extraction. J. Mach. Learn. Res. **4**, 177–210 (2003)
8. Choi, F.Y.: Advances in domain independent linear text segmentation. In: Proceedings of the 1st North American Chapter of the Association for Computational Linguistics Conference (NAACL), pp. 26–33. Association for Computational Linguistics, Seattle, Washington, USA, April 2000
9. Cohen, W., Ravikumar, P., Fienberg, S.: A comparison of string metrics for matching names and records. In: KDD Workshop on Data Cleaning and Object Consolidation, vol. 3, pp. 73–78 (2003)
10. Eisenstein, J., Barzilay, R.: Bayesian unsupervised topic segmentation. In: Proceedings of EMNLP 2008, pp. 334–343 (2008)
11. Elving, W.J.L., Westhoff, J.J.C., Meeusen, K., Schoonderbeek, J.-W.: The war for talent? The relevance of employer branding in job advertisements for becoming an employer of choice. J. Brand Manag. **20**(5), 355–373 (2013)
12. Ferret, O., Grau, B., Masson, N.: Thematic segmentation of texts: two methods for two kinds of texts. In: Proceedings of the 36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics, vol. 1, pp. 392–396. Association for Computational Linguistics (1998)
13. Hall, M., et al.: The Weka data mining software: an update. ACM SIGKDD Explor. Newsl. **11**(1), 10–18 (2009)
14. Hearst, M.A.: Texttiling: segmenting text into multi-paragraph subtopic passages. Comput. Linguist. **23**(1), 33–64 (1997)
15. Kessler, R., Torres-Moreno, J.M., El-Bèze, M.: E-gen: automatic job offer processing system for human resources. In: Proceedings of the 6th Mexican International Conference on Artificial Intelligence, pp. 985–995 (2007)
16. Kuhn, P., Mansour, H.: Is internet job search still ineffective? Econ. J. **124**(581), 1213–1233 (2014)
17. Lafferty, J.D., McCallum, A., Pereira, F.C.N.: Conditional random fields: probabilistic models for segmenting and labeling sequence data. In: Proceedings of the 18th International Conference on Machine Learning, San Francisco, USA, pp. 282–289 (2001)

18. Loth, R., Battistelli, D., Chaumartin, F.-R., de Mazancourt, H., Minel, J.-L., Vinckx, A.: Linguistic information extraction for job ads (SIRE project). In: Adaptivity, Personalization and Fusion of Heterogeneous Information, Paris, France, pp. 222–224 (2010)
19. McCallum, A., Li, W.: Early results for named entity recognition with conditional random fields, feature induction and web-enhanced lexicons. In: Proceedings of HLT-NAACL, vol. 4, pp. 188–191. ACM (2003)
20. Misra, H., Yvon, F., Jose, J.M., Cappe, O.: Text segmentation via topic modeling: an analytical study. In: Proceedings of the 18th ACM Conference on Information and Knowledge Management, pp. 1553–1556. ACM (2009)
21. Müller, O., et al.: Towards a typology of business process management professionals: identifying patterns of competences through latent semantic analysis. Enterp. Inf. Syst. **10**(1), 50–80 (2016)
22. Peng, F., McCallum, A.: Information extraction from research papers using conditional random fields. Inf.. Process. Manag. **42**(4), 963–979 (2006)
23. Pevzner, L., Hearst, M.A.: A critique and improvement of an evaluation metric for text segmentation. Comput. Linguist. **28**(1), 19–36 (2002)
24. Ponte, J.M., Croft, W.B.: Text segmentation by topic. In: Peters, C., Thanos, C. (eds.) ECDL 1997. LNCS, vol. 1324, pp. 113–125. Springer, Heidelberg (1997). https://doi.org/10.1007/BFb0026725
25. Reynar, J.C.: Statistical models for topic segmentation. In: Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics, pp. 357–364. ACM, Maryland, USA (1999)
26. Riedl, M., Biemann, C.: Topictiling: a text segmentation algorithm based on LDA. In: Proceedings of ACL 2012 Student Research Workshop, pp. 37–42. ACM (2012)
27. Scaiano, M., Inkpen, D.: Getting more from segmentation evaluation. In: Proceedings of HLT-NAACL, pp. 362–366. ACM, Montreal, Canada (2012)
28. Sha, F., Pereira, F.: Shallow parsing with conditional random fields. In: Proceedings of HLT-NAACL, pp. 134–141. ACM, Edmonton, Canada (2003)
29. Simon, A.-R., Gravier, G., Sébillot, P.: Leveraging lexical cohesion and disruption for topic segmentation. In: International Conference on Empirical Methods in Natural Language Processing, EMNLP 2013, pp. 1314–1324 (2013)
30. Stamatatos, E., et al.: Clustering by authorship within and across documents. In: Working Notes Papers of the CLEF 2016 Evaluation Labs, September 2016
31. Tschuggnall, M., et al.: Overview of the author identification task at PAN-2017: style breach detection and author clustering. In: Working Notes Papers of the CLEF: Evaluation Labs, p. 2017. CLEF and CEUR-WS.org, September 2017
32. Utiyama, M., Isahara, H.: A statistical model for domain-independent text segmentation. In: Proceedings of the 39th Annual Meeting on Association for Computational Linguistics, pp. 499–506 (2001)
33. Winkler, W.E.: The state of record linkage and current research problems. In Statistical Research Division, US Census Bureau. Citeseer (1999)
34. Zavrel, J., Berck, P., Lavrijssen, W.: Information extraction by text classification: corpus mining for features. In: Proceedings of the Workshop Information Extraction Meets Corpus Linguistics (2000)