



Universiteit
Leiden
The Netherlands

Algorithmic tools for data-oriented law enforcement

Cocx, T.K.

Citation

Cocx, T. K. (2009, December 2). *Algorithmic tools for data-oriented law enforcement*. Retrieved from <https://hdl.handle.net/1887/14450>

Version: Corrected Publisher's Version

License: [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#)

Downloaded from: <https://hdl.handle.net/1887/14450>

Note: To cite this publication please use the final published version (if applicable).

Chapter 5

Enhancing the Automated Analysis of Criminal Careers

Four enhancements have been devised, both on a semantic and efficiency level, that improve existing methods of automated criminal career analysis described in Chapter 4. A new distance measure was introduced that more closely resembles the reality of policing. Instead of the previously suggested, more rigid, comparison of career changes over time we propose an alignment of these careers. We employ a faster and better method to cluster them into a two-dimensional plane. This visualization can ultimately be used to predict the unfolding of a starting career with high confidence, by means of a new visual extrapolation technique. This chapter discusses the applicability of these new methods in the field and shows some preliminary results.

5.1 Introduction

In Chapter 4 (cf. [7]), we described research attempting to gain insights into the concept of criminal careers: the criminal activities that a single individual exhibits throughout his or her life. The resulting tool analyzed the criminal record database in described in Appendix B. This method mainly addressed the extraction of the four important factors (see Section 5.2) in criminal careers and established an overview picture on the different existing types of criminal careers by using a stepwise or year-based approach, with the ultimate objective to prepare the data for prediction of a new offender's career. The approach centered on formalizing an individual's *criminal profile* per year, representing an entire career as a string of these calculated profiles. Chapter 4 identified some difficulties in comparing these strings and provided solutions to them. The method, however, suffered from time-complexity issues and a comparison mechanism that was largely rigid and ad hoc.

In this chapter we describe a number of techniques were specifically developed for the enhancement of the above mentioned analysis but are also widely applicable in other

areas. We explain how these methods can be used to improve the existing paradigm by replacing the individual ad hoc methodologies and how combining them will reduce the number of steps needed to reach better results. We supplement the existing approach by adding a preliminary prediction engine that employs the properties of the already realized visualization. Experiments performed with this setup are also discussed.

The main contribution of this research lies in the novel combination of a number of separately created technologies, all very well suited for law enforcement purposes, to pursue the common goal of early warning systems that allow police agencies to prevent the unfolding of possibly dangerous criminal careers.

Four enhancements have been devised. In Section 5.3 a new distance measure is introduced that more closely resembles the reality of policing. Instead of previous, more rigid, comparison of career changes over time we propose an alignment of these careers in Section 5.4. In Section 5.5 we employ a faster and better method to cluster them into a two-dimensional plane. This visualization can ultimately be used to predict the unfolding of a starting career with high confidence, by means of a new visual extrapolation technique (see Section 5.6). Section 5.7 shows experiments, and Section 5.8 concludes.

5.2 Background and Overview

In Chapter 4 the four factors were extracted from the criminal record database, recorded as numbers and treated as such. However, it appears to be beneficial to treat a collection of crimes in a single year as a *multiset*, which then describes severity, nature and frequency inherently. A multiset or *bag*, is a collection where each element can occur more than once. The set of all distinct elements in that multiset is called its *underlying set*. Figure 5.1 describes the relation between the two and shows how we employ it to represent a criminal's activities in a single year.

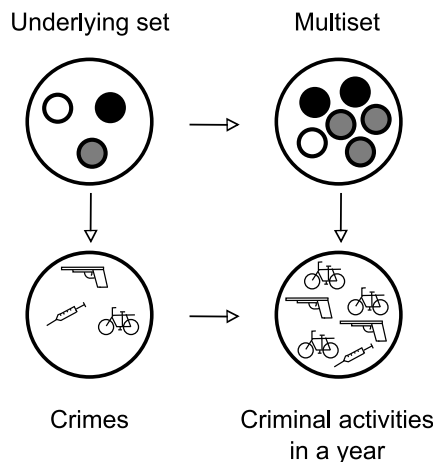


Figure 5.1: A multiset representation of a criminal profile in a single year

The multiset representation offers advantages, most notably the availability of standard approaches to compare multisets and calculate distances between them. Kusters and Laros [18] devised a distance function for multisets that generalizes well-known distance measures like the *Jaccard distance* [17]. This metric contains a customizable function f that can be adapted to fit specific knowledge domains. It also allows for the incorporation of weights for an element, e.g., element x counts twice as much as element y . We employ this metric for calculating the difference between two crime-multisets (see Section 5.3) by choosing a specific f .

Instead of the former method of a strict number-wise comparison between years (comparing the first year of criminal a with the first year of criminal b , the second year of a with the second year of b , etc.), with the possibility of stretching or shrinking careers (cf. Section 4.5, we propose a novel *alignment* of the mentioned multisets. This method strives for an optimal automated matching of years, using the distance measure described above, assigning penalties for every mutation needed, which enables a police analyst to better cope with situations like captivity, forced inactivity or unpenalized behavior. Section 5.4 elaborates on this.

Visualization and clustering methods used before yielded good results, mainly by incorporating direct input from police specialist. It was however computationally complex, taking a very long time to process the entire dataset. Within the scope of the research into criminal careers we developed a method that improved standard *push-and-pull* algorithms but eliminated the need for human input, while retaining the former strength of its output [19]. We explain this in Section 5.5.

Earlier results provided some important information for law enforcement personnel, but the ultimate goal, predicting if certain offenders are likely to become (heavy) career criminals, was not realized or elaborated upon. Further investigation of this possibility led to the need of a good two-dimensional visual *extrapolation* system, described in [12]. The conceivement of this technique paved the way for possible development of early warning systems, that we explore in Section 5.6.

5.3 A Distance Measure for Profiling

The core of our approach consists of a new distance measure. This metric is a special case of the *generic metric for multisets* as described in [18]:

$$d_f(X, Y) = \frac{\sum_i f(x_i, y_i)}{|S(X) \cup S(Y)|}$$

This metric calculates the distance between two finite multisets X and Y , where $S(A) \subseteq \{1, 2, \dots, n\}$ is the underlying set of multiset A and a_i is the number of occurrences of element i in multiset A . Here $f : \mathbb{R}_{\geq 0} \times \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{\geq 0}$ is a fixed function with finite supremum M and the following properties:

$$\begin{aligned} f(x, y) &= f(y, x) && \text{for all } x, y \in \mathbb{R}_{\geq 0} \\ f(x, x) &= 0 && \text{for all } x \in \mathbb{R}_{\geq 0} \\ f(x, 0) &\geq M/2 && \text{for all } x \in \mathbb{R}_{> 0} \\ f(x, y) &\leq f(x, z) + f(z, y) && \text{for all } x, y, z \in \mathbb{R}_{\geq 0} \end{aligned}$$

These properties ensure that d_f is a valid metric [18].

Naturally, all results depend largely upon choosing the right *defining function* f for a certain knowledge domain. In the law enforcement area, it is important to realize that the relative difference between occurrences of a crime is more important than the difference alone. This is because the distance between an innocent person and a one-time offender should be much larger than for example the distance between two career criminals committing 9 and 10 crimes of the same kind respectively, thus ensuring that $f(0, 1) \gg f(9, 10)$ (the two arguments of f are the numbers of respective crimes of a given category, for two persons).

A good candidate for this function, that was developed in cooperation with police data analysts, seems to be the function

$$f_{crime}(x, y) = \frac{|x - y|}{(x + 1)(y + 1)}$$

for integer arguments x and y , both ≥ 0 . This function complies with the above mentioned characteristic of criminals and yields a valid metric.

It is obviously important to still be able to incorporate *crime severity* and *nature*, next to crime frequency, into this equation. This is possible through the addition of *weights* to the generic metric described above. A suggestion was made by Kusters and Laros [18] for accomplishing this: each item i gets an assigned integer weight ≥ 1 and is multiplied by this weight in every multiset, simulating the situation where all weights were equal but there were simply more of these items in each set. An impression of how this affects the calculation of distances between criminal activities is given in Figure 5.2.

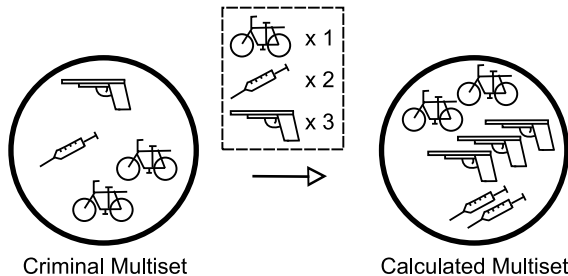


Figure 5.2: Adding weights to criminal activities

By using this specifically tailored distance measure with weighing possibilities, we are able to calculate distances between criminals per single time frame. This distance can serve as basis to discover distances between careers, that can be seen as strings of these time frames.

5.4 Alignment of Criminal Careers

Since the temporal ordering of crimes is of high importance in the definition of a criminal career, a tool occupied with its analysis should provide a way to calculate distances

between these ordered series based upon distances between their elements. Sequence alignment could be a valuable tool in this process.

In contrast with the strict year-by-year comparison used in Chapter 4, the alignment paradigm [21] tries to find a best match between two ordered series, allowing a small number of *edits* to be made in such an ordering: insertions, deletions and other simple manipulations. The penalty for a deletion or insertion is governed by a so-called *gap-penalty function*. Alignment is an algorithm typically used within the area of computational biology. Famous instances include those of Needleman-Wunsch [23] and Smith-Waterman [27]. Each alignment is based upon a valid metric for elements (a year of accumulated crimes in this case), as for example the metric discussed in Section 5.3. Figure 5.3 describes a typical alignment situation, showing two such edits. The treatment of gaps, however, may somewhat differ from biological applications; also note that empty multisets (an “innocent” year) have a special relation with the gap penalty. Next to these differences, careers may or may not be “complete”; some of them are still unfolding.

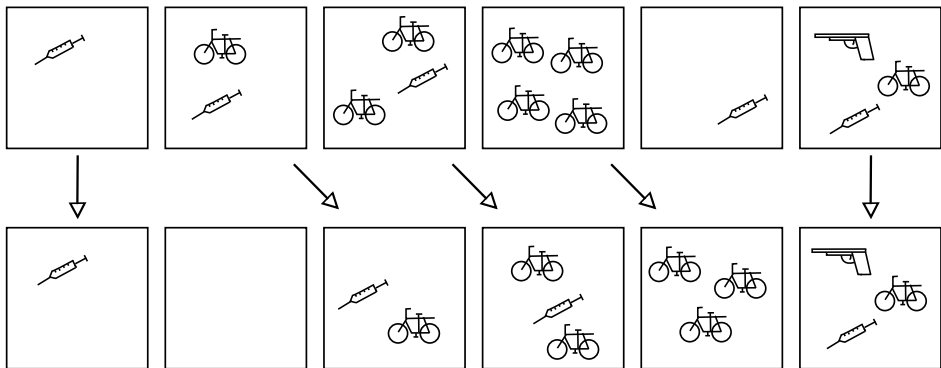


Figure 5.3: Two criminal careers whose similarity is revealed by alignment

One of the rationales for using alignment instead of strict sequence matching is the frequent occurrence of gaps or stretches in criminal data. One could think of, for example, prison time, forced inactivity, judicial struggles consuming time between an offense and its penalization and, most importantly, the frequent occurrence of unpenalized criminal activity due to undiscovered crime. Treating this naturally occurring plethora of phenomena without any descriptive value on one’s intentions in the criminal area as non-existent or not important will result in a deformed image of compared criminal careers. Therefore the usage of an alignment mechanism is strongly favored over a strict comparison on a semantic level. Another reason for this might be the occurrence of randomly appearing year changes within one’s career, e.g., two crimes occurring either in June and July, or in January and December, are in essence the same, although strict separation between years will place the latter in two different time frames. Although data is often gathered in such broad time frames, leading to the inevitable occurrence of this problem, the alignment paradigm could potentially be used to reduce the subsequent effects.

Careful consideration needs to be given to the gap-penalty function in the field of

crime analysis. While each gap in standard (biological) alignment procedures represents a constant penalty, a gap in a criminal career can be considered less or more important based upon the amount of time passed within this gap. We therefore assign a time stamp $t(a)$ to each element a in the criminal career. The gap-penalty is then computed by applying the gap-penalty function to the different $\Delta t(u, i) = t(u_{i+1}) - t(u_i)$ involved, where u_i is the i^{th} element in u , an ordered, time stamped series. The known algorithms to compute distances have to be adapted, causing an increase in time ($O(n^2) \rightarrow O(n^3)$) and space complexity ($O(n) \rightarrow O(n^2)$), where n is the length of the series.

Using this alignment we compare all possible couples of criminal careers and construct a standard distance matrix of the results.

5.5 Toward a New Visualization Method

In the original setup, a *push-and-pull algorithm* [5] was used to create a (sub)optimal two-dimensional visualization of the produced distances matrix from a randomized starting point. Initially, points get random positions in the unit square of the plane; after which they are repeatedly pulled together or pushed apart, depending on the relation between current and desired distance. This method, a variant to Multi-Dimensional Scaling, relied upon domain knowledge provided by the analyst using the tool. This domain input realized a significant gain in results compared with simpler, unsupervised algorithms of the same kind [20].

The complexity of the previous algorithm, however, prohibited a time efficient usage of the tool, costing near days on standard equipment to analyze the target database and thus posing a serious bottleneck in the original setup that needed to be overcome, preserving the property that elements can be easily added and traced.

One of the major problems of the simpler variants of the push-and-pull algorithm was the fact that a large number of randomly positioned individuals tended to create subgroups that pushed a certain item harder away from its desired position in the plane, than comparable items pulled this point toward that position (see Figure 5.4, standard).

Especially in the case of criminal careers, where analysis tends to result in vastly different, very large clusters, this effect appears to present. Addressing this issue could therefore lead to a representative visualization and good response times, eliminating the need for human input. Kusters and Laros [19] suggested the usage of a *torus* rather than the standard bounded flat surface. Within such a torus, all boundaries are identified with their respective opposite boundaries, enabling the “movement” of an individual from one end of the surface to the opposite end, thus overcoming the above mentioned problem (Figure 5.4, torus).

When using this torus to construct a two-dimensional visualization of our distance matrix, time complexity was reduced to the level of simple push-and-pull algorithms, but the visualization is expected to be of the same quality as the image produced by the supervised algorithm. Other MDS techniques are also potential candidates for adoption, if they adhere to the demand of incremental addition of single items.

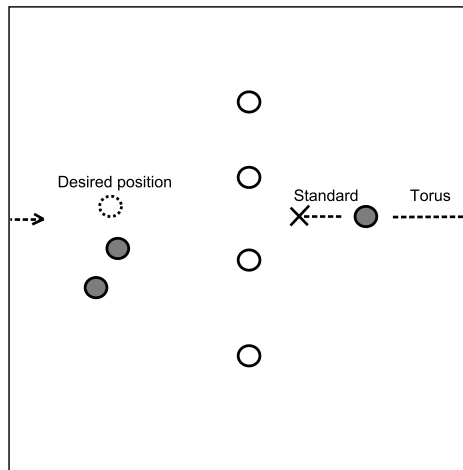


Figure 5.4: Advantage of torus visualization versus standard visualization

5.6 Prediction of Unfolding Careers

As was mentioned in Chapter 4, a lot of information is inherently present in the visualization of mutual distances. In Chapter 7, [12] we propose to utilize the power of a two-dimensional visualization for temporal extrapolation purposes. This algorithm plots the first, already known, time frames of criminal activity of a certain person in the correct place (compared to other fully grown careers already displayed in the image) and extrapolates these points in the plane to automatically discover careers that are likely to become very similar in the future. A system employing this should be able to provide police operatives with a warning when such a starting career can easily develop into that of a heavy career criminal. More on the possibilities concerning prediction can be found in Chapters 7 and 8, where an actual attempt is made at the prediction of criminal careers based upon this data.

5.7 Experimental Results

We tested our new approach on the criminal record database (cf. Appendix B), containing approximately one million offenders and the crimes they committed. Our tool analyzed the entire set of criminals and presented the user with a clear two-dimensional visualization of criminal careers as can be seen in Figure 5.5. During this process, a multiset string of offenses was added to each offender, where each bucket contained crimes committed within a certain time frame.

The image in Figure 5.5 gives an impression of the center of the torus output produced by our tool when analyzing the before mentioned database. This image shows the identification that could easily be coupled to the appearing clusters after examination of

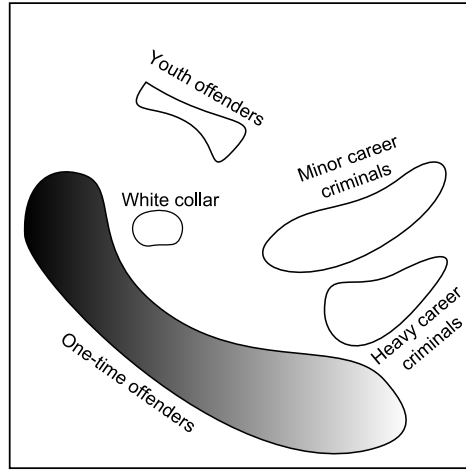


Figure 5.5: Impression of clustered groups of criminal careers in the visualization

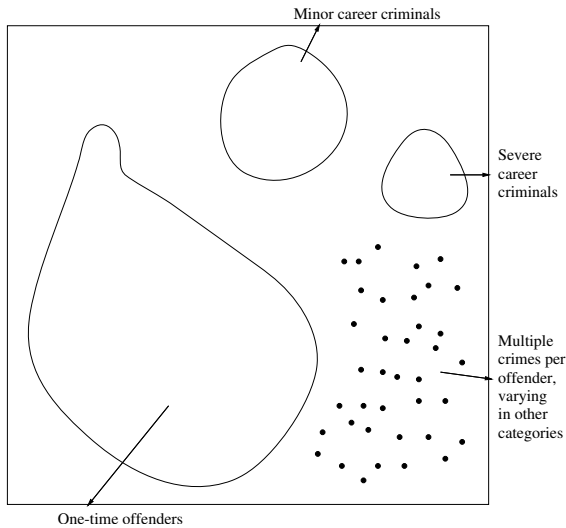
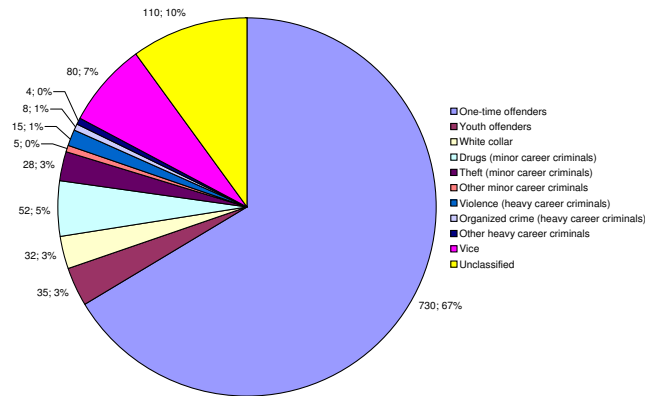


Figure 5.6: Results from Chapter 4

its members. In total 10 different clusters were identified using a standard k-means clustering method, of which some were left out of Figure 5.5, because they were too small or because a combination of clusters provided a visualization that was easier to interpret. Below, in Figure 5.7 they are specified in more detail, accompanied by the distribution of all offenders from the database over these classes, where amounts are specified in thousands.

These figures appear to be describing reality better than the visualization resulting



<i>Career Class</i>	<i>Percentage</i>
One-time	66%
Youth Offenders	3%
White Collar	3%
Drugs (Minor Career)	5%
Theft (Minor Career)	3%
Other Minor Career	1%
Violence (Heavy Career)	1%
Organized Crime (Heavy Career)	1%
Other Heavy Career	1%
Vice	7%
Unclassified	10%

Figure 5.7: Proportion of offenders in each class of criminal career

from Chapter 4, [7]. Just like the image constructed by the previous approach (Figure 5.6) the new image features a large “cloud” of one-time offenders. However a clear distinction can now be noticed between the left, dark, part of that cloud, which represent minor crimes, and the right, lighter side of the cluster that contains the heavier felonies. Next to this, the group of miscellaneous offenders was split into two, more specific, groups; the white border criminals and youth offenders. The remaining unclustered individuals were left out of the image for clarity reasons. Analysis of this group can however also reveal interesting results, answering why these individuals do not belong to one of the larger clusters and what kind of crimes these people commit. Next to the better results provided by this approach, far better computation times were realized that outperform the previous method by a factor 100.

Getting more insights into the possible existence of subgroups in any cluster remains a desirable functionality of this approach. Future research will focus on getting this improvement realized (cf. Section 5.8). Next to that, it might be of interest to find common subcareers in the complete list of careers, to find out if there are any subcareers that occur

often. An even more important search would be to find careers that are common in one class of offenders and are not in others, in a way representing the class they are common in. These possibilities are discussed in Chapter 6.

5.8 Conclusion and Future Directions

In this chapter we demonstrated the applicability of some newly developed methods in the field of automated criminal career analysis. Each of the enhancements provided a significant gain in computational complexity or improved the analysis on a semantic level. An integral part of the new setup consists of the multiset metric that was specifically tuned toward the comparison of criminal activities. It leaned strongly upon knowledge provided by domain experts. This distance measure was used to search for optimal alignments between criminal careers. This edit distance for sorted or temporal multiset data improved greatly upon the original year-by-year comparison, dealing with problematic situations like captivity or forced inactivity. Using this alignment a distance matrix can be constructed and clustered on a flat surface using a new unsupervised method, yielding comparable results with the previous setup, but largely reducing the time complexity. The power of this visual representation was employed to extrapolate a starting career within the plane of the visualization, predicting its unfolding with some accuracy. The end report consists of a visual two-dimensional visualization of the results and a prototype of an early warning system that predicts newly developing criminal careers, which is ready to be used by police experts.

This chapter describes improvements on a methodology for the analysis of criminal careers introduced in Chapter 4. Building upon a series of successfully introduced algorithms, this novel approach is subject to a high number of parameters and leans heavily on the way these techniques work together. It is to be expected that the current setup can also be improved by tweaking its parameters and elaborating on the internal cooperation between phases. Future work will therefore focus on a larger experimental setup to find and verify optimal settings and retrieve results that are even more descriptive and usable. We also hope to equip our tool with a sub-cluster detection algorithm to provide even better insights into the comparability of criminal careers.

It may be of interest to set fuzzy borders between the different years. Crimes within months ending or beginning such a time frame can be (partly) assigned to the next or previous year respectively as well, thus eliminating the problems arising with strict coherence to the change of calendar year.

Special attention will be directed toward the predictability of criminal careers, and the eventual suitability of this approach for early warning systems running at police headquarters throughout the districts. Incorporation of this new tool in a data mining framework for automatic police analysis of their data sources is also a future topic of interest.

As this chapter describes an extension to the material described in Chapter 4, the same privacy and judicial constraints apply to the matter discussed here. They are discussed in detail in Appendix A.