How applicable are attribute-based approaches for human-centered ranking creation?

Clara-Maria Barth, Jenny Schmid, Ibrahim Al-Hazwani, Madhav Sachdeva, Lena Cibulski, Jürgen Bernard



 PII:
 S0097-8493(23)00059-6

 DOI:
 https://doi.org/10.1016/j.cag.2023.05.004

 Reference:
 CAG 3681

To appear in: *Computers & Graphics*

Received date : 15 October 2022 Revised date : 23 April 2023 Accepted date : 10 May 2023

Please cite this article as: C.-M. Barth, J. Schmid, I. Al-Hazwani et al., How applicable are attribute-based approaches for human-centered ranking creation?. *Computers & Graphics* (2023), doi: https://doi.org/10.1016/j.cag.2023.05.004.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2023 Published by Elsevier Ltd.

PDF of Manuscript, including all figures, tables etc. - must not contain any author details

Click here to view linked References ±

Computers & Graphics (2023)



How Applicable are Attribute-Based Approaches for Human-Centered Ranking Creation?

Clara-Maria Barth^{a,*}, Jenny Schmid^a, Ibrahim Al-Hazwani^{a,d}, Madhav Sachdeva^a, Lena Cibulski^{b,c}, Jürgen Bernard^{a,d}

^aUniversity of Zurich, Zurich, 8006, Switzerland

^bFraunhofer IGD, Darmstadt, 64283, Germany

^cTechnical University of Darmstadt, Darmstadt, 64289, Germany

^dDigital Society Initiative, Zurich, 8006, Switzerland

ARTICLE INFO

Article history: Received April 23, 2023

Keywords: Item Ranking, Attribute-Based Ranking, Human-Centered Ranking Creation, Human Factors, Visual Analytics, Experimental Study

ABSTRACT

Item rankings are useful when a decision needs to be made, especially if there are multiple attributes to be considered. However, existing tools do not support both categorical and numerical attributes, require programming expertise for expressing preferences on attributes, do not offer instant feedback, lack flexibility in expressing various types of user preferences, or do not support all mandatory steps in the ranking-creation workflow. In this work, we present RankASco: a human-centered visual analytics approach that supports the interactive and visual creation of rankings. The iterative design process resulted in different visual interfaces that enable users to formalize their preferences based on a taxonomy of attribute scoring functions. RankASco enables broad user groups to a) select attributes of interest, b) express preferences on attributes through interactively tailored scoring functions, and c) analyze and refine item ranking results. We validate RankASco in a user study with 24 participants in comparison to a general purpose tool. We report on commonalities and differences with respect to usefulness and usability and ultimately present three personas that characterize common user behavior in ranking-creation. On the human factors side, we have also identified a series of interesting behavioral variables that have an influence on the task performance and may shape the design of human-centered ranking solutions in the future.

© 2023 Elsevier B.V. All rights reserved.

11

12

13

1. Introduction

In everyday life, people constantly face the challenge of finding the best item in an item set: whether it is about picking the nicest hotel for a holiday trip, the next movie to watch, the most promising stock to buy, or the perfect flat to rent. Item sets typically contain large numbers of items to choose from, each of which is defined across multiple attributes representing different criteria to be considered carefully. Such a multi-attribute choice [1] is not an easy task, especially for non-experts. An obvious optimal solution does generally not exist, and "best"

*Corresponding author: e-mail: clara-maria.barth@uzh.ch (Clara-Maria Barth) highly depends on the decision-maker's personal preferences. Also, the task complexity heavily depends on the dataset size and the number of relevant attributes, both subject to growth.

A strategy to identify items of choice in large item sets is 14 the creation of item rankings. A striking benefit of rankings is 15 the inherent order they provide to items, enabling users to eas-16 ily find most preferred items at the top. In turn, least preferred 17 items for a decision-making scenario are situated at the bottom 18 of the ranking. We focus on human-centered approaches for the 19 creation of item rankings, leveraging individual preferences of 20 users as a profound basis to express multiple criteria to opti-21 mize for. Traditionally, many people relied on pen and paper or 22 general purpose spreadsheet tools to formalize and create item 23 rankings. With the digital transformation, people can make use 24 of more sophisticated computational support to ease the cre-25

Preprint Submitted for review / Computers & Graphics (2023)

ation of item rankings. Still, interactively engaging with the

2 creation and refinement of item rankings is desirable for every-

one: not only for domain experts or users with programming

⁴ expertise but also for non-experts.

Strategies for interactive ranking creation are two-fold. Itembased approaches allow users to express feedback about the perceived order and relevance of data items [2, 3, 4, 5]. Users can directly interact with items of interest, make item com-8 parisons, and adjust the ranks of items, e.g., in spreadsheets. Attribute-based approaches allow users to express preferences 10 on attributes. Algorithms then transform these preferences into 11 attribute scores, combine the attribute scores according to some 12 weighting, and produce a ranking as the direct result of order-13 ing items by their overall scores [6, 7]. With the proposal of 14 Attribute Scoring Functions (ASFs) [8], we have presented the 15 formal underpinning to define user preferences on attributes. 16 In this work, we provide an extension to RankASco [9], a 17 visual analytics tool around ASFs to create attribute-based item 18

rankings. We focus on the attribute-based creation of item rank-19 ings for two reasons. First, we believe that its scalability is 20 mainly agnostic to the number of items, making it more appli-21 cable for large item sets. Second, we assume that it is easier 22 for users to express preferences for individual attributes than 23 between items as a whole. A pioneer visual analytics approach 24 25 for multi-attribute ranking is LineUp [6]. It offers a visual interface that allows users to map attribute values to preference 26 scores, even if LineUp does not offer full flexibility regarding 27 types of user preferences. 28

The reflection on the body of related work on ranking cre-29 ation in general revealed five shortcomings. First, existing so-30 lutions do not yet offer the flexibility users may require to in-31 tuitively express their preferences regarding attribute values. In 32 33 specific, we identify a lack of tools that support both categorical and numerical ASFs. Second, most existing tools require 34 programming when it comes to ASF creation. These tools can 35 only be steered by math experts or computer scientists, but not by non-experts. Third, the black-box nature of the program-37 ming paradigm does not offer instant feedback about the dis-38 tribution of attribute values (data), how a created ASF behaves 39 (model), or how interactive refinements by users affect the pro-40 cess (user). Fourth, hardly any tool supports all mandatory 41 steps in the ranking workflow: creation, refinement, and us-42 age [8]. Finally, not much is known about the users of ranking 43 creation tools. In particular, a deeper understanding about com-44 mon ranking creation behaviors could help the design and de-45 velopment of (human-centered) ranking systems in the future. 46 In this context, the evaluation of approaches for the characteri-47 zation of user groups, such as personas, could be useful to better 48 understand user needs when creating rankings. 49 To this end, we revisit and extend RankASco [9]. Our con-50

51 tributions are as follows:

 The presentation of RankASco, an attribute-based visual analytics approach that accepts user preferences to create rankings for large item sets. RankASco is the result of a two-year

research project, with two workshop paper publications [8, 9] forming the baseline for this extended version. We build

forming the baseline for this extended version. We build
 upon RankASco with additional visual interfaces, refined de-

sign choices, and more descriptive details.

• The validation of RankASco in a user study with 24 partici-

pants. The study evaluates RankASco in comparison to Excel as a representative of a general purpose tool. We decided to recruit non-experts with low familiarity in using programmatic solutions to solve multi-criteria ranking problems.

• The presentation of three personas, characterizing common user behaviors in ranking-creation: (1) Peter, the perfectionist, (2) Eva, the explorer, and (3) Pippa, the pragmatist.

By providing visual interfaces for all eight types of attribute 67 scoring functions, our approach is the first that allows users to 68 express a large variety of attribute-based preferences, for cate-69 gorical and numerical attributes alike. In an iterative design pro-70 cess, we have developed RankASco to make the task of rank-71 ing creation accessible to a broad range of users. As a result 72 of careful design, development, and validation, RankASco pro-73 vides a framework that supports multi-criteria decision-making 74 for the general public. With the identification of three personas 75 for item ranking, we hope to guide the design of future human-76 centered ranking solutions. 77

2. Related Work

Ranking creation typically relies on algorithmic models that 79 leverage data characteristics to infer an item order. We ex-80 tend this principle towards human-centered creation of rank-81 ings, which encourages users to interactively engage with the 82 underlying data and express preferences on items or attributes. 83 We structure related works along our main contributions: the 84 general role of human preference expression (Section 2.1), the 85 human-centered creation of item rankings (Section 2.2), and the 86 evaluation of interactive approaches for ranking creation and 87 personalization (Section 2.3). 88

78

89

106

2.1. Expression of User Preferences

Providing the users with the ability to input their preferences 90 is a crucial aspect of human-centered design in various fields, 91 such as recommender systems, visual analytics, and human-92 computer interaction. There are two main approaches to gather 93 user preferences: implicit feedback and explicit feedback. Im-94 plicit feedback is based on collecting information about the 95 users' preferences by watching their natural interaction with 96 the systems, e.g., number of clicks or time spent on a page. 97 Explicit feedback requires the users to explicitly express feed-98 back, e.g., by selecting and marking documents and providing 99 ratings for specific items. The main advantage of implicit feed-100 back is that there is no cost for the user to provide feedback. 101 However, it is generally thought that the implicit strategy tends 102 to be less accurate than explicit feedback [10]. For a detailed 103 comparison between these two approaches, please refer to ex-104 isting studies [11, 12]. 105

2.2. Human-Centered Ranking Creation

Ranking creation in real-world settings is mostly performed by third-party platforms, thus leaving users only with the resulting ranking. Most web shops, movie streaming services, and online browsing follow this line of approach. Algorithmic support for ranking creation often involves recommender systems [13, 14] or other types of machine learning methods [15, 16, 17]. We exclude this branch of approaches, as it does

Preprint Submitted for review/Computers & Graphics (2023)

not allow users to explicitly create rankings by themselves,
 thus not following the human-centered principle. In fact, some
 third-party approaches enable users to *personalize* existing
 rankings, but they do not allow for initial ranking creation. In
 contrast, human-centered ranking creation offers a high degree
 of human control [18], where they can apply preferences either
 to items or attributes. This distinction structures our reflection
 on related works.

Item-based approaches allow users to explicitly express feedback about items and their perceived order to arrive at a 10 personalized ranking. TasteWeights [2] enables users to itera-11 tively adjust item preferences using slider widgets. While users 12 can directly observe how their modification of preferences 13 affects the ranking, the approach cannot assign negative item 14 preference scores to indicate disfavor. RanKit [4] exploits the 15 users' knowledge at an item level by providing a user-friendly 16 interface for users to manually rank known items. 17 As a 18 beneficial side effect, the authors identified an increase in user trust towards the resulting ranking, which stems from real-time 19 visual feedback on user's interactions. Finally, Podium [19] 20 21 is a multi-attribute approach that enables users to drag items across the ranking to reflect the perceived relative relevance 22 23 of items. Podium then infers the parametrization of a ranking SVM model to match these preferences. To complement the 24 computational support, users can also change the weights 25 of attributes contributing to the item ranking. Off-the-shelf 26 spreadsheet approaches such as Microsoft Excel, Google Sheet, 27 and Apple Numbers can be seen as item-based approaches, en-28 abling users to perform analysis tasks like filtering and sorting. 29 Manageable task complexity depends on the user's level of 30 expertise: if users are required to solve a complex ranking task, 31 some considerable scripting skills will be required. 32

Attribute-based approaches allow users to explicitly express 33 preferences regarding specific attributes and attribute values of 34 items. In previous work [8], we studied different approaches 35 that can be used for transforming attribute values into scores, 36 ranging from merely theoretical approaches [20, 21] to visual 37 interactive approaches [22, 23, 24]. The resulting taxonomy of 38 eight types of attribute scoring functions serves as a baseline in 39 this work to study human-centered ranking creation based on 40 attribute preferences. A pioneer work for attribute-based rank-41 42 ing creation is LineUp [6], an interactive technique designed to 43 create, visualize, and explore rankings of items based on a set of heterogeneous attributes. LineUp enables users to formalize 44 functions that map attribute values to scores, either through a 45 46 programming interface or through visual interfaces. The visual approach supports the formalization of linear and compound 47 linear (e.g., a roof-function) preferences. However, no inter-48 active visual support is provided for discontinuous functions 49 or categorical attributes. MyMovieMixer [25] is an interactive 50 movie recommender system. Users can select filter criteria and 51 apply linear item preferences by using a slider widget. The 52 authors report that users perceived to be more in control of the 53 ranking results by expressing their preferences explicitly. How-54 ever, MyMovieMixer does not support non-linear preferences. 55 56 WeightLifter [26] is an interactive visualization that allows users to explore the relationship between attribute weights and 57 ranking results, thus increasing the transparency of the ranking 58 model. Users can simultaneously explore up to 10 attributes. 59

However, trade-offs between more than two attributes require 60 attribute grouping to weigh them via sliders, making it difficult 61 for users to precisely express their preferences. Moreover, 62 WeightLifter assumes that attribute values do not require trans-63 formation beyond normalization to be considered as attribute 64 scores. RankViz [5] is a visualization framework that enables 65 users to compare two rankings and see how each attribute has 66 contributed to the items' ranking positions. Its major downside 67 is that it requires users to have some knowledge about ranking 68 algorithms, thus shifting the focus from a more personalized 69 ranking towards a more interpretable ranking model. uRank [7] 70 is an interactive approach for understanding, refining, and 71 reorganizing document items on-the-fly as information needs 72 evolve. Specifically, it enhances predictability through docu-73 ment hint previews, which serve two purposes: allowing users 74 to control the ranking by choosing keywords and supporting 75 understanding by means of a transparent visual representation 76 of scores. To summarize, while promising attribute-based 77 approaches exist, none of the reviewed approaches supports 78 users in expressing all types of desirable preferences [8]. To be 79 able to study commonalities and differences among item-based 80 and attribute-based approaches, we present an extension of 81 RankASco to be used in our proposed experimental study. 82

2.3. Evaluation of Human-Centered Ranking Creation

Approaches for the human-centered creation of rankings are commonly evaluated with usage scenarios [27] and qualitative experiments [28], such as user studies.

Usage scenarios report on how a proposed approach could 87 be used, highlighting the strengths of the approach in solving 88 a specific task. For example, the evaluation of RanKit [4] 89 employed a usage scenario to clearly illustrate the steps from se-90 lecting a dataset selection to showing how user feedback is used 91 to improve the ranking. Similarly, Podium [3] leverages a usage 92 scenario to showcase how the approach can be used to identify 93 the most important features of the user's favorite football team. 94

Qualitative experiments are used to observe and collect feed-95 back on how users interact with an approach in a real-world 96 setting [28]. Item-based approaches have been evaluated by 97 recruiting a number of participants, including both experts and 98 non-experts. The experiments use pre- and post-questionnaires 99 to understand more about how users solve assigned tasks. 100 Attribute-based approaches have been mostly evaluated 101 with expert users, as in the case of WeightLifter [26] and 102 RankViz [5]. One reason may be that the tasks that users aimed 103 to solve have been mostly technical to date. For example, to 104 evaluate RankViz [5], knowledge about ranking algorithms was 105 required to fully understand also the non-visual mechanics. 106

Hardly any studies have compared item-based approaches 107 with attribute-based approaches. So far, the visualization 108 community does not offer reflections on commonalities and dif-109 ferences of the two types of approaches, and designers of visu-110 alization approaches for the creation of ranking algorithms rely 111 on their experiences when it comes to task abstractions, require-112 ment engineering, and iterative visualization and interaction 113 design. A pioneer evaluation approach has been taken by Gratzl 114 et al. with LineUp [6]. The authors conducted a pre-study with 115 just experts using item-based approaches like Microsoft Excel 116 or Tableau, and a post-study with expert and non-expert users 117

3

83

84

85

using the proposed attribute-based approach. The studies highlighted that novice users were faster in solving the task using LineUp compared to experts using Microsoft Excel or Tableau. Our experiment goes beyond this scope, as we analyze across-subject item agreement, task completion time, and de-5 rive personas as a reflection of our behavioral observations. The usage of personas to characterize user behavior is a wellknown method in HCI research and practice, such as system design [29], product design [30], and marketing [31]. A persona represents a user group's unique collection of behavior 10 patterns, objectives, and talents as a realistic character to make 11 them more actionable and understandable [32]. 12

3. Scoring Functions for Attribute-Based Ranking 13

The attribute-based ranking approach leverages user pref-14 erences regarding attribute characteristics. Expressions auto-15 16 matically have an effect on all items, regardless of the dataset size. To rank items based on attribute preferences, attribute val-17 ues must be transformed into numerical values that represent 18 the preference scores of users. We call this process attribute 19 scoring. For example, users preferring fast cars might favor 20 21 high HP attribute values, while penalizing low HP attribute values.Ultimately, all attribute-based preference scores can be used 22 and combined to create the overall item ranking. 23

To perform the mapping from attribute values to preference 24 scores, we build upon Attribute Scoring Functions (ASFs) [8], 25 serving as one of our two baseline workshop publications that 26 we extend in this work. We briefly echo the essentials of ASFs, 27 which are described and discussed in the baseline work in de-28 tail. In short, ASFs are mappings of data attributes that: 29

• transform the input values to numerical output scores, 30

• have a *polarity* for the output score domain, and 31

• have a valence for the output scores. 32

Data Transformation Each ASF covers the entire input do-33 main of an attribute. This ensures that each attribute value can 34 be mapped to an output score. In addition, any attribute value 35 must be mapped to exactly one output score to ensure the valid-36 ity of the data transformation and to prevent ambiguity. 37

Polarity The output score domain of an ASF has a pre-38 defined range. Similar to normalization, these pre-defined 39 40 ranges allow for comparable preference scores across attributes. Value ranges can either be uni-polar (e.g., ranging from 0 to +1) 41 or bi-polar (e.g., ranging from -1 to +1). Having a uni-polar 42 range for the output allows users to express how much they like 43 attribute values, while a bi-polar range also allows users to ex-44 press how much they dislike certain attribute values. 45

Valence Output scores of ASFs carry valence information, 46 which implies that each output score has semantic meaning. On 47 the one hand, higher scores always represent higher preferences 48 of users compared to lower scores. On the other hand, extreme 49 scores (possibly caused by extreme input values) automatically 50 imply stronger preference values. 51

We differentiate between categorical and numerical ASFs to 52 53 explicitly account for the different characteristics of categorical and numerical data attributes. In total, we identified and 54 described eight different types of ASFs in our taxonomy pre-55 sented in the baseline work, shown in Figure 1. Three ASFs 56



Fig. 1. Taxonomy of eight types of ASFs, used as a functional baseline [8].

are applicable to categorical attributes and five are applicable 57 to numerical attributes. For the sake of self-explainability, we 58 briefly re-iterate the eight types. 59

3.1. Categorical Attributes

Categorical ASFs can be used for the transformation of 6 categorical attributes to preference scores. There are three 62 different types of categorical ASFs, which are explained in 63 the following sections: Score Assignment, Equidistant, and 64 Non-Equidistant [8]. 65



Score Assignment. Score Assignment ASFs are the simplest type of categorical ASFs. They work based 67 on absolute preferences, where users directly assign 68 an absolute preference score to each category, in the 69 notion of an explicit quantification [33] of categorical values. 70

60

76

77

78

87

88

89

90

This ASF type can be used for assigning exact preference scores 71 to all categories. These scores are absolute, meaning that users 72 can assign preference scores without comparing different cate-73 gories. A real-world example includes the assignment of scores 74 to different holiday destination cities. 75



Equidistant. Equidistant ASFs can be used for the assignment of relative preferences to categories. With this ASF type, users can create an order of all cate-

gories and assign preference scores to the categories, 79 according to their position in the overall order. The equidis-80 tant ASF distributes the score values equally across the value 81 domain. This can be useful if users know about the preferred 82 order of categories, but cannot express how much they prefer a 83 certain category over another. A real-world example includes 84 ordering of different colors for furniture, where users are sure 85 about the order of colors. 86



Non-Equidistant. Non-Equidistant ASF extend the precision of Equidistant ASFs. They also work based on relative preferences but allow for non-equidistant value score distributions between ordered categories.

Especially when several categories of an attribute appear to be 91 similar, non-equidistant ASFs enable users to also assign simi-92 lar preference scores. For a movie example, the non-equidistant 93 ASF can be used for the ordering of movie genres where users 94 may like very few genres. 95

3.2. Numerical Attributes

4. Abstractions

Numerical ASFs can be used for the mapping of continuous numerical attribute values to preference scores using numerical functions. There are five types of numerical ASFs: Two-Point Linear, Two-Point Non-Linear, Multi-Point Continuous, Multi-5 Point Discontinuous, and Quantile Based [8], all described in the following sections.



12

Two-Point Linear. Two-Point Linear ASFs are a simple type of numerical ASFs and can be used for expressing linear preferences where attribute values at the end of the range (on the top or bottom end) can be favored. This ASF type is suited for non-complex preferences.

Examples from the mathematical domain include the min-max 13 14 or max-min normalization. Real-world examples include the preference for cheapest prices for mobile phone subscriptions. 15

Two-Point Non-Linear. Two-Point Non-Linear ASFs 16 17 consist of two points at the start and end of the input range and a line segment in between but, contrary 18 to the Two-Point Linear ASF, can reflect non-linear 19 preferences. This allows users to steer the skewness of the 20 underlying attribute value distribution, enabling users to create 21 ASFs that are similar to, e.g., logarithmic functions or the 22 square root norm. A real-world example is the logarithmic 23 preferences for TV screen sizes, where above a certain point an 24 increase in screen size is only a marginal improvement. 25

Multi-Point Continuous. Multi-Point Continuous 26 27 ASFs expand the design space of ASFs considerably through the addition of additional points within the 28 input value range. Therefore, they allow the creation 29 30 of more complex and even compound functions. Multi-point Continuous ASFs can reflect sophisticated user preferences 31 that are not monotonically increasing or decreasing, such as 32 preferences for middle values (i.e., roof-like functions) or ramp 33 functions. A real-world example is a preference for middle-34 priced shoes, since they often have the best price-quality ratio. 35

37 38 39

Multi-Point Discontinuous. Multi-Point Discontinuous ASFs introduce the concept of mathematical discontinuities to the ASF design space. In Multi-Point Discontinuous ASFs not all points must be con-

nected, allowing the creation of functions with gaps in the out-40 put domain. A mathematical example of this behavior is a stair 41 function. Real-world examples include the preference for either 42 old-timer cars or the latest car models at the same time (with 43 low preferences for middle-aged cars). 44

Quantile Based. Quantile Based ASFs are different 45 from the Two-Point and Multi-Point ASF types in that 46 they apply statistical quantile normalization to the at-47 tribute values. In contrast to value-based functions, 48 49 the order of values determines the output scores of distribution, similar to the notion of a rolling pin for baking. This ASF type 50 allows users to flatten narrow value distributions, and limit the 51 impact of outliers in the dataset. 52

We briefly characterize the main steps of the workflow when 54 performing attribute-based creations of item rankings, before 55 we describe the rationales that motivated the design of our 56 visual analytics approach. The driving principle was the strin-57 gent support for users to express their subjective preferences on 58 attributes, following the goal to create a human-based data an-59 alytics solution. The ranking creation workflow is inspired by 60 the work of Wang et al. [34], Kuhlman et al. [35], Gratzl et al. 61 [6], and Cheng et al. [36]. Since our approach is based on user 62 preferences, preferences are the basis of the attribute selection 63 rather than automated selection as in Wang et al. [34]. Overall, 64 we have identified three principal phases in the workflow to 65 create a human-centered ranking, as Figure 2 illustrates. 66

- 1. Attribute Overview and Selection: Users should first gain an overview of attributes and select interesting attributes.
- 2. Creation of ASFs: For each selected attribute, users can create an ASF such that their preference for certain attribute values can influence the ranking.
- 3. Ranking Analysis: The ranking is presented to users, enabling the analysis of the validity of the computed ranks.

We articulate seven requirements to visual analytics approaches for the human-centered interactive visual creation of item rankings. These requirements are based on the problem statement, related work on multi-criteria decision-making [6, 26, 37, 5, 36], experiences gained through previous work [8], and by echoing human-centered visual analytics principles:

- R1: Attribute Overview: Providing an overview of attributes, their value distributions, and dependencies between attributes to support the informed selection of attributes.
- R2: User Preferences: Accounting for individual user preferences, creating various ASF types should be supported.
- R3: Instant Feedback: Assessing the effect of changed ASFs on underlying data distribution values instantly should be supported for validation and refinement purposes.
- R4: Straight-Forward ASF Creation: Opening attribute scoring to a diverse spectrum of users should be supported.
- R5: Ranking Overview: Analyzing the ranking results, including influencing scores, should be possible.
- R6: Attribute Weighting: Approaches should allow defining and refining the importance of attributes through weights, achieving human-centered rankings.
- R7: Ranking-Data Comparison: Assessing how the ranking relates to the underlying item distribution should be possible for users.

5. RankASco - Human-Centered Attribute-Based Ranking

We present a visual analytics approach to support users 99 in the interactive creation of human-centered item rankings. 100 RankASco (short for Ranking based on Attribute Scorings) is 101 an attribute-based approach that takes users preferences into 102 account to calculate an item ranking, even for large datasets. 103 We present a refined and extended version of RankASco, as an 104 extension of the original workshop paper publication [9]. An 105 overview of the three main views of RankASco can be seen in 106

67

68

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

Overview of	all attribute	Attribute Sco	bring Function	Resu	lt Overv	view	Download JS0	Download Op	ptions Downlo	oad Weights	Download CSV	back		
Categorical Attribute	es			PCA of Input	t distribution	b	aseRent	condition	yearConstructed	d living	Space PCA of	Output distribution		
noParkSpaces	condition	noRooms	regio	Control Liner	There are Diverse with a straining value.	-1 0.5			0.85	0.71	0.28	0.	40 0.5	0.5 1
Add	Add	Add	Add			Rank	ID	Color	Item Score	Item Score Symbol	baseRent	condition	yearConstructed	livingSpace
Numerical Attributes				°402 221 1011	100 203	2	113596488		1.182		0.925	0.522	0.687	-0.178
baseRent 🔳	yearConstructed	livingSpace	serviceCharge			3	106922601		1.078		0.489	0.522	0.687	0.243
				Attributo Sco	coring Eulection	4	115370074		1.047		0.450	0.522	0.778	0.185
	100			in a second second	Output der before	5	82393787		0.792		0.078	0.522	0.999	0.185
932 2170	1930 2020	45 200	13 532.75			6	110307393		0.909		0.346	0.522	0.778	0.064
Missing Values: 0	Missing Values: 0	Missing Values: 0	Missing Values: 0			7	106864949		0.925		0.346	0.522	0.563	0.253
Add	Add	Add Correlation between Cates	Add			8	114547944		0.932		0.450	0.522	0.778	-0.101
Correacon between numerical attributes (rearson correacion) Correacioni detween categorical attributes (chi Square)				Score Assignment Doublater1		8	114547944		0.932		0.450	0.522	0.778	-0.101
byingdigana	• 5	pardonOrtBalcory-		 Non-Equidistant wird_condition 		10	115106022		0.716		-0.002	0.522	0.778	0.320
yearContracted baseRent serviceCharge EvingSpace	x yearCaristructed basekent			-veljest -inkrbined 1 pesperimed) 2 2	ulind_time_see	1 2	3 4 5 6	7 8 5	10 Next Last					
	(1	(2	2)					(3)						

Fig. 2. Overview of the RankASco visual analytics workflow. Users can (1) gain an overview of multiple categorical and numerical attributes and underlying correlations between attributes, (2) create attribute scorings for relevant attributes based on their preferences by using interactive visual interfaces, and (3) configure attribute weights, analyze and refine ranking results, and make informed multi-criteria decisions.

Figure 2, in line with the three main phases of the workflow

proposed in Section 4. RankASco is publicly available https: 2 //rankasco-ivda.ifi.uzh.ch/, with more implementation

details in the supplemental material.

4

6

5.1. Phase 1: Attribute Overview and Selection 5

The first phase of the interactive workflow for the creation 6 of item rankings consists of the identification of a meaningful set of attributes that are relevant to the users' preferences. The 8 attribute overview and the correlation overview allow users to 9 make an informed selection on a set of attributes (R1). The 10 attribute overview interface in RankASco shows all existing 11 attributes for a given item set, as shown in Figure 2 (left). 12 For categorical attributes, all categories and their counts are 13 shown in bar charts. For numerical attributes, histograms show 14 15 the distribution of numerical values. In addition, RankASco also reveals the number of missing values for each attribute. 16 The handling of missing values is crucial to calculate an item 17 score for each item. LineUp [6] handles missing values by 18 calculating the mean or median of an attribute; we use an 19 approach where users can define a score for missing values 20 explicitly. The handling of missing values is different for cat-21 egorical and numerical attributes, as described in the respective 22 sections. When users select a set of interesting attributes, there 23 likely exist correlations between attributes. To account for this 24 important decision-making criterion, the extended version of 25 RankASco now offers a correlation overview for categorical 26 attributes and for numerical attributes alike, shown in Figure 2 27 28 step 1 (bottom). Categorical correlations are calculated based on the Chi-squared test [38], while the Pearson correlation 29 coefficient [39] is used for numerical correlations. 30

5.2. Phase 2: Creation of Attribute Scoring Functions 31

After the identification of a set of relevant attributes, users 32 can create an ASF for each attribute to use for the calculation 33

of the item ranking. The selected types of ASFs are based on 34 the eight different types of ASF that we identified in a baseline 35 work [8]. RankASco supports this stage by providing eight dif-36 ferent interactive visual interfaces for the creation of the eight 37 different types of ASFs, which is the core of the baseline pub-38 lication [9]. With the eight visual interfaces, a broad spectrum 39 of mental models of users can be addressed (R2): Some ASF-40 creation interfaces are simple and straightforward, while other 41 variants are more complex and highly customizable. To guide 42 users in the selection and the creation of an ASF, visual finger-43 prints explain the functional behavior of the ASFs and respec-44 tive interfaces. This helps users find the best ASF type for their 45 preferences and the underlying attribute data. 46

The design of all eight ASF interfaces follows the same prin-47 ciples: Input values (the attribute values) are shown on top left 48 in the ASF creation view, output values (the output scores) are 49 shown on top right, next to the input values as can be seen in 50 Figure 4 (left). This eases the comparison between the char-51 acteristics of the input and output value distribution, and thus 52 the effects of the ASF on the data attribute. The actual ASF-53 creation interface is always shown below the two distribution 54 charts and differs for all eight types of ASFs. The iterative pro-55 cess particularly focused on the design of interfaces that are 56 easy to use (R4). Direct manipulation and linking of views 57 update the output value distribution in real-time whenever the 58 ASF is modified (R3). Design and implementation details of 59 the eight interfaces for ASF creation are as follows. 60

5.2.1. Categorical Attributes

Three interfaces allow users to create ASFs with preference 62 scores for categorical attributes. All three interfaces share the 63 same strategy to support missing value treatment: missing val-64 ues of categorical attributes always form an additional category 65 that can be considered by users for scoring purposes. 66

61

67

The Score Assignment ASF is based on absolute preferences.

Attribute Scoring Function		
Investigation in the second		Two-Point NonLinear
• Destilation Kinds	Esslingen_Kreis	1,
 Reutlingen_Kreis 	Ludwigshurg Kreis	
Schwarzwald_Baar_Kreis	• Luungabuig_ticia	0
	Stuttgart	
 Waldshut_Kreis 		-1
ess preferred	0 more preferred	

Fig. 3. Two different interfaces for ASF creation. The categorical Non-Equidistant ASF is used to, e.g., express strong preferences for three regions in an apartment-hunting situation (left). A numerical Two-Point Non-Linear ASF shows users preferences for low service charges for the apartment of choice.

It allows users to assign a numerical preference score to each
category. In RankASco, this ASF type is represented through
numerical input fields (one for each category in a categorical
attribute) where users can directly assign preference scores between -1 and +1. looseness=-1 To ease the usage, users can
also start with a pre-defined neutral value for all categories and
assign preference scores for a subset of categories only. This
feature overcomes the need for setting a score for every cate-

⁹ gory, even if irrelevant. This is especially efficient for categori-10 cal attributes of high cardinality.

The Equidistant and Non-Equidistant ASF types are based on 11 relative preferences. The interface of both ASF types are two-12 13 dimensional, where categories are shown along the y-axis and preference scores are shown along the x-axis. Users can adjust 14 the position of each category by horizontal dragging interaction, 15 from the left (less preferred) to the right (more preferred). The 16 difference between the two ASFs is the placement strategy of 17 categories along the x-axis: for the Equidistant ASF, categories 18 are positioned along discrete equidistant positions, to guarantee 19 equal spacing between categories. In contrast, with the Non-20 *Equidistant* ASF, users have the ability to position categories 21 continuously along the x-axis, allowing for non-equal spacing 22 between categories. Figure 3 (left) shows a Non-Equidistant 23 ASF that represents users preferences for certain regions. A de-24 tailed example for the Score Assignment and Equidistant ASFs 25 can be found in the supplemental material. 26

27 5.2.2. Numerical Attributes

The interfaces for the four value-based numerical ASF 28 types (all numerical ASFs except Quantile based) use a two-29 dimensional coordinate system. Attribute values are shown on 30 the x-axis and preference scores are shown on the y-axis. This 31 design choice is based on mathematical functions f(x) = y and 32 how they are visualized in 2D. The interfaces for each ASF type 33 initialize a default function that can be adjusted with draggable 34 points, e.g., to steer the slope and curvature of the line segments 35 between points. The user-created function determines how in-36 put values are transformed into preference scores. An example 37 38 of each of the four numerical ASFs can be found in the supple-39 mental material, including an enlarged figure. The Two-Point Linear ASF consists of one linear line seg-40

41 ment spanning across the entire input value domain. Two points
 42 at the very left and very right of the x-axis can be vertically ad-

justed, to change the slope of the ASF. An example of a *Two-Point Linear* ASF can be found in the supplemental material.

The *Two-Point Non-Linear* ASF type expands this concept by allowing for a non-linear line segment between the two points. The curvature of this line segment can also be steered through an additional point, a so-called control point, based on the mathematical concept of Bézier curves [40]. Figure 3 (right) shows a Two-Point Non-Linear ASF that shows a non-linear (logarithmic) preference for service charge values.

Multi-Point ASF types have more than two points and also more line segments, respectively. Thus, they allow more flexibility in function design. The mode of operation is the same as for the *Two-Point* ASFs: The line segments and their curvature can be steered through draggable points, as shown in Figure 4. For the *Multi-Point Continuous* ASF, all lines are always connected to each other, resulting in a continuous mathematical function. *Multi-Point Discontinuous* ASFs, on the other hand, introduce mathematical discontinuities between line segments to create gaps in the output domain.

The *Quantile Based* ASF works based on statistical quantile normalization, which is applied to the order of the attribute values instead of the actual values. One insight we had in the design process was to allow steering the degree to which an input value distribution shall be subject to quantile normalization. The interface now offers a slider that lets users steer the degree of quantile normalization that is applied to the data, ranging from 0% (no quantile normalization at all) to 100% (full quantile normalization applied).

5.3. Phase 3: Ranking Analysis

The final phase of the ranking creation workflow is the analysis, validation, and possible refinement of the created ranking. Based on the set of created ASFs and a user-steerable weight for each attribute, an overall item score is calculated for each item. This score is a weighted sum of all the attribute preference scores multiplied by their attribute weight (more details are given below). The final item ranking then results from ordering all item scores in decreasing order.

The design of the ranking interface is inspired by list-based item visualizations, typically utilized in interfaces for search results [41, 42, 43], and the output of recommender systems [2, 44, 45]. The ranking result is split into multiple pages, which allows users to either only look into the top items or, if interested, 84

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78



Fig. 4. One of five visual interfaces for the creation of numerical ASFs. Here, a Multi-Point Continuous function is used to represent the user's preferences for base rent prices around €1300 (drag-and-drop interface on the right). On the left, two histograms show the distribution of input values and output scores. This instant feedback also helps to achieve balanced scores, compared to the left-skewed input values.

also check items ranked in the middle by using the pagination. This visualization allows to only print the top items first and, if requested, load additional items to handle large item sets better. To account for details that help explain item ranks, we extend the list-based idea to a tabular layout, as shown in Figure 5 (right). This table contains all items (rows) as well as columns 6 for the attribute scores involved in the human-centered ranking process (R5). Users can steer the weights for all attributes with 8 sliders, as shown in Figure 5 (left). These adjustable weight sliders allow users to assign preferences of importance to the 10 different attributes and are initially set to 0.5. Modifying one of 11 the attribute weights results in a re-calculation of all item scores 12 13 in real-time and an update of the item ranking (R6). The weighted score for each item is supported with a visual 14 cue that eases the comparison of different items, as shown in 15 Figure 5. Additional scatter plots support the comparison of the 16 input and output data characteristics (R7): The input data (Fig-17 ure 5 left) consists of a dimensionality-reduced version of the 18 one-hot encoded original dataset. The output data (Figure 5 19 right) consists of a dimensionality-reduced version of all at-20 tribute scores for each item. Every item also has a unique ID 21 and a color that is determined in a similarity-preserving way. 22 Colors are assigned based on a 2D color map [46] on either the 23 input or output distribution of all items, as shown in the scatter 24 plots. Linking items between one of the two scatter plots and 25 the ranking result facilitates the comparison of top-ranked items 26 and their distribution originating from the input or output data 27 distributions. To finish the interactive ranking creation work-28 flow, users can export the ranking (CSV or JSON format) and 29 use it for downstream analyses, as shown in Figure 5 (right). 30

6. Usage Scenarios: Apartment-Hunting

8

We introduce two usage scenarios for multi-criteria ranking problems with a dataset about apartments in the south-west of Germany, publicly available at Kaggle [47]. We will be talking about the Fischer family: Hugo Fischer, husband of Barbara Fischer and father of two girls. Hugo is a fictive non-expert who is looking for a new apartment for his family. Overall, the Fischer family has ten preferences on apartments with different priority, pertaining to ten different attributes. This scenario will recur in our experimental study; details on exact preferences of the Fischer family are described in the supplemental material.

6.1. Apartment-Hunting with the general purpose tool

Hugo takes an item-based approach, using a spreadsheet tool such as Microsoft Excel, to organize the apartment items in a preferred way. He decides for the spreadsheet tool because he owns a software license for the tool anyway and has been using Excel occasionally during the last years. The tabular format gives him control over the items (rows), while always having the lookup of attributes (columns).

42

43

44

45

46

47

48

49

First, he starts with gaining an overview of the dataset. 50 Vertical scrolling helps him traverse all items, and he realizes 51 that the number of available apartment items is large. Hor-52 izontal scrolling lets Hugo identify preferred attributes and 53 delete irrelevant attributes to reduce task complexity. Next, 54 he uses a filter to reduce the dataset size: for the region, he 55 excludes regions other than "Stuttgart", "Ludwigsburg Kreis", 56 and "Esslingen Kreis", which reduces the dataset size by 57 roughly 80%. Next, Hugo uses a filter in the notion of a 58 dynamic query [48] to remove *base rent* values below €1000 59 and above €1600, as the allocated budget of the family is about 60 €1300. As a third operation, Hugo selects the rooms column 61 and filters out room sizes outside the values 3.5, 4, 4.5, and 62 5. He then sorts the service cost charges from least to most, 63 as low additional costs are important for the family. After 64 these operations, Hugo notices that only 26 flats remain for 65 decision-making. Possibly, Hugo's filter criteria have been too 66 stringent for one or the other attribute, such that valuable items 67 may have disappeared; he therefore relaxes the filter criteria 68 a bit. Also, Hugo is not yet happy with the ranking order of 69

Preprint Submitted for review/Computers & Graphics (2023)

Result Overview		Download JSON	Download Option	5 Dowr	nload Weights	Download CSV	Back					
Please select whether the color map sho	uld be calculated based on the input	or the output.		Rank	ID	Color	Item Score 👌	Item Score Symbol	baseRent	condition	yearConstructed	livingSpace
 Input 				1	111161515		1.182		0.507	0.522	0.778	0.402
ං Output				2	113596488		1.280		0.925	0.522	0.687	-0.178
PCA of Input distribution	baseRent	condition	PCA of Output distribution	3	106922601		1.078		0.489	0.522	0.687	0.243
· · · · · · · ·	0	0		4	115370074		1.047		0.450	0.522	0.778	0.185
	vearConstructed	livingSpace	0.5	5	82393787		0.792		0.078	0.522	0.999	0.185
-1 0.5 1 1.5 2	0	0 -1 0.40		1 2	3 4 5 6	7 8	9 10 Next Last					

Fig. 5. Ranking result overview with color-coding across views. The interface on the left shows the data spaces of the input and the output scores (scatter plots) and allows steering attribute weights. The ranking result interface on the right shows details on attribute scores, with similarity-preserving colors.

items. Hugo decides to sort by the *base rent* attribute due to its
high importance and discovers that the attribute has a bipolar
nature, meaning that the best apartments around €1300 are
not at the top. With some scripting effort, he fixes the problem
and, with some additional scripting, manages to sort items by

more than one attribute at the same time. Finally, to take all ten

⁷ preferences on apartments into account, he starts with changing

⁸ the order of items manually to arrive at a final ranking. Hugo

shows the result to Barbara, and together they determine which

¹⁰ of the top flats they want to visit as a family.

11 6.2. Apartment-Hunting with RankASco

We demonstrate the usefulness of RankASco for the multicriteria ranking problem of the Fischer family. We will accompany Hugo's workflow until he has created ASFs for the four preferences of highest priority.

Hugo starts using RankASco by analyzing all attributes and 16 attribute value distributions shown in Figure 2 (left). He is par-17 18 ticularly interested in the region, base rent, number of rooms, and *service charge*; so he starts with the *region* attribute. He 19 creates the Non-Equidistant ASF shown in Figure 3 (left) rep-20 21 resenting his strong preference for the three regions "Stuttgart", "Ludwigsburg Kreis", and "Esslingen Kreis" (in that order). 22 Next, Hugo picks the base rent attribute and creates the Multi-23 Point Continuous ASF, depicted in Figure 4 (right). The roof-24 like function punishes apartments that are too cheap. Begin-25 ning with €1000, apartments turn positive, with a maximum 26 at €1300. Even larger prices for rent turn into negative scores 27 28 at €1600. Then, Hugo defines his preference on 4-room flats (with 3 to 5 rooms also deemed acceptable), using a Score As-29 signment ASF. Next, Hugo chooses the service charge attribute 30 with a preference for service charges as low as possible, rep-31 resented with a Non-Linear ASF shown in Figure 3 (right). 32 The non-linear nature of the function returns positive scores for 33 many of the low values of service charge, but decreases steeply 34 for very high values. This is an example of how Hugo can ex-35 ploit the bipolar support for scores given with RankASco (po-36 larity characteristics), 37 After creating the four ASFs, Hugo proceeds to the ranking 38 overview and starts with refining the attribute weights per 39 attribute, as shown in Figure 5 (left). Given his preference 40 scores and weights per attribute, RankASco automatically 41

⁴² provides the resulting item ranking (Figure 5 (right)). From ⁴³ here, Hugo's remaining process is three-fold. First, Hugo can

refine the four created ASFs, if the analysis of the ranking

⁴⁵ result reveals aspects that can be improved. Second, he uses

RankASco's export functionality to show the preliminary list of top candidates to his wife Barbara. Third, he continues with adding the missing six preferences to arrive at the final ranking according to the Fischer's preferences, as a start for the informed visit of quasi-optimal apartments. 50

7. User Study

With the proposal of a visual analytics approach for the 52 human-based creation of item rankings, we widen the band-53 width of existing approaches in a still loosely populated design 54 space. Interesting questions emerge regarding the evaluation 55 of RankASco, but also with respect to human factors involved 56 in the ranking creation realm. For that purpose, we conducted 57 an experimental study with two distinct parts: the observation 58 and analysis of user performance, and the observation of and 59 reflection on user behavior. The first main goal of our experi-60 ment was to compare RankASco with a general purpose tool, 61 similar to the user study of LineUp [6] with Excel and Tableau. 62 We crosscut this performance analysis with our second goal: to 63 observe and identify user behaviors among study participants 64 to ultimately derive personas. In the study, data collection in-65 cluded quantitative and qualitative data by taking participants' 66 task completion time, determining across-subject item agree-67 ment in the top 20 ranking results, recording behavioral obser-68 vations, and conducting informal interviews. We first describe 69 the research questions and the experiment design, before we 70 provide details on the results of the two study parts in Section 8. 71

7.1. Research Questions

The two main goals can be broken down into four research questions as follows:

- *RQ*₁: Can a stringently attribute-based ranking approach compete with the general purpose tool in terms of efficiency?
- RQ_2 : Does the number of items have an impact on the performance of the attribute-based ranking creation in comparison to the general purpose tool?
- RQ_3 : How do users behave in the three different phases of the ranking workflow?
- *RQ*₄: Is it possible to derive personas from observed user behavior in both RankASco and the general purpose tool?

 RQ_1 and RQ_2 are related to the quantitative assessment of user performance (Part 1), while RQ_3 and RQ_4 are related to behavioral observations of users (Part 2).

51

72

73

74

75

76

77

78

80

81

82

Preprint Submitted for review / Computers & Graphics (2023)

7.5. Task Description

7.2. Experiment Factor: RankASco and General Purpose Tool

We use RankASco as a representative of a multi-criteria 2 attribute-based ranking tool for the interactive creation of item rankings. The decision for RankASco is based on its completeness in the support of categorical and numerical attributes 5 through interactive visual interfaces to support eight types of ASFs and its stringent design for large user groups, including non-experts. In contrast to, e.g., LineUp [6], RankASco is the only attribute-based ranking approach that entirely works without the need for coding. It also supports numerical and 10 categorical attributes. To allow for a fair comparison between 11 RankASco and the general purpose tool, the correlation plots in 12 RankASco were disabled for the user study. 13

We used a general purpose tool as a representative of the dif-14 ferent options with which users can create item rankings in their 15 everyday live. Those include filtering by many and sorting by 16 one attribute. We aimed for an approach for the creation of item 17 rankings that should neither require expert knowledge nor pro-18 gramming skills. An overview of items and attributes should 19 be provided, such that users can make informed decisions on 20 the ordering of items. Finally, users should be able to express 21 22 preferences through direct manipulation. Similar to Gratzl et al. [6], we used Excel due to its popularity for the targeted user 23 24 population.

7.3. Experiment Factor: Dataset Size 25

The underlying dataset forms the basis for a second experi-26 ment factor: the dataset size. We utilized the Kaggle "Apart-27 ment Rental Offers in Germany" [47] dataset for the exper-28 iment. Overall, the dataset contains 268,850 apartment list-29 ings (items) with a total of 49 attributes. After the exclusion 30 of binary, range, redundant, ambiguous, and task-irrelevant at-31 tributes, we chose six categorical and four numerical attributes 32 for this study. In an upstream process, we made data quality 33 checks and eliminated items that contained null values, missing 34 attributes, or implausible values (cf. supplemental material). 35 To study user performance with respect to the dataset size, 36

we control the number of items as one experiment factor. From 37 the remaining items of sufficient quality, we randomly selected 38 39 500 items. To arrive at different experiment conditions, we used these items to create three subsets: 500 items, 300 items, and 40 100 items. 41

7.4. Participant Description 42

We recruited 24 participants (14 female) at the university, 43 aged between 22 and 31 (M = 26, SD = 3.09). A prereq-44 uisite for the experiment was a basic command of Excel and 45 the ability to understand and speak the offered experiment lan-46 guages (EN & DE). Human subjects research approval from 47 the faculty's ethics board was obtained prior to the study. Par-48 ticipants who completed the study received a gift card worth 49 \$/€30 as compensation. Prior to the study, we asked partici-50 pants about their knowledge in Excel (M = 3.29, SD = 0.94), 51 data science (M = 3.29, SD = 0.84), and multivariate data 52 analysis (M = 2.79, SD = 1.00), using a 5-point Likert scale 53 54 (high signifies very good knowledge). Additionally, we asked whether the participants had prior experience creating rankings 55 for decision-making problems (29% yes) or have already solved 56 a ranking problem programmatically (25% yes). 57

The task for all participants was to create a ranking for a 59 given set of items and the tool at hand. To facilitate the com-60 parability of results, we controlled the preferences that partici-61 pants would have to follow in the experiment, i.e., we introduce 62 the truth of the ranking scenario upfront. We designed a nar-63 rative evaluation [49], where participants assumed the role of a 64 real estate agent who is aiming at identifying the top 20 apart-65 ment items, based on the preferences of the Fischer family, the 66 clients of the real estate agent (cf. the Usage Scenario in Sec-67 tion 6 and the supplemental material). We designed the prefer-68 ences of the Fischer family based on two main goals. First, the 69 preferences should include a healthy mix of attribute character-70 istics, with preferences ranging from simple (e.g., "the higher, 71 the better") to complex (e.g., mathematical discontinuities like 72 "preferred if apartment is built before 1900 or as new as pos-73 sible"). Second, the preferences should all be plausible for a 74 family with two kids. Overall, ten preferences needed to be 75 considered, each for a different attribute. We created a tabular 76 description of the preference attributes, sorted by their impor-77 tance. 78

7.6. Dependent and Independent Variables

Independent variables are the type of ranking approach (us-80 ing the general purpose tool or the RankASco) and the dataset 81 size (100, 300, or 500 items). The crosscut of the two variables 82 leads to six experiment conditions. Every participant was as-83 signed to only one dataset size level, i.e., eight participants were 84 tasked with 100 items, etc. In contrast, every participant was 85 asked to perform the ranking task on both types of ranking ap-86 proaches to maximize the comparability across approaches. To 87 avoid the learning effects and effects of fatigue, we randomized 88 the tool to start with between participants for all three dataset 89 sizes. Dependent variables are the task completion time and the 90 across-subject item agreement of the top 20 ranking results. 9

79

92

93

94

95

96

97

98

99

7.7. Study Procedure

We carried out a pilot study in advance to make sure that task and study design were understandable, robust, and feasible. The study procedure included four steps: (1) introduction, (2) training, (3) ranking creation, and (4) questionnaire. We introduced participants to their task by providing a narrative and dataset description, to ease the lookup of preferences of the Fischer family whenever needed (cf. supplemental material). Then, we conducted an introduction session to the approaches used 100 so that participants could always familiarize themselves. Par-101 ticipants were trained by walking them through the interfaces 102 of the tools. Using the movies dataset we explained each ASF 103 with a usage example. 104

In the core of the study, participants solved the ranking task 105 with the two approaches. In parallel, we conducted an observa-106 tional study to also assess the user behavior. We measured the 107 participants' task completion time without prior announcement, 108 to avoid time pressure on the participants' side. To assess 109 across-subject item agreement, we collected the final top 20 110 items after participants completed the ranking task. Finally, 111 we conducted a qualitative interview utilizing a 5-point Likert 112 scale rating to assess and compare participants' perceived 113 confidence in their ranking results (cf. supplemental material). 114

Preprint Submitted for review/Computers & Graphics (2023)



Fig. 6. (A) Comparison of task completion time of RankASco versus general purpose tool. (B) Relative difference of task completion time, depending on the number of items. Values > 0 indicate that using RankASco is faster, whereas values < 0 indicate that using the general purpose tool (Excel) is preferable. (C) Task completion time for the subsets of 100, 300, and 500 items. (D) The heterogeneity of the top 20 ranked items across participants withing the 100 items dataset (for the 300, 500 items see the supplemental materials). Here, 1 indicates that no other participant agreed with that item being in the top 20 ranks (high heterogeneity), whereas 8 indicates that a particular item was in the top 20 items of all participants (low heterogeneity).

Interview questions also included the users' experience with

7.8. Data Analysis

Part 1: Performance Analysis (RQ_1 , RQ_2). We analyzed the 5 performance measures for a) the comparison of the two item ranking approaches (RQ_1) , b) the comparison of the three dataset sizes (RQ_2) , and c) the cross-cut of both experiment factors (2x3 conditions) (RQ_2). To perform this data analysis strategy, we used a two-fold approach. First, we used visual 10 representations for (a-c) to assess effects visually (see Figure 6). 11 Second, we applied statistical tests to identify considerable 12 or even significant differences between the conditions with 13 respect to the dependent variables. The test portfolio included 14 a paired two-sample t-test [50] and Wilcoxon Rank-Sum 15 Test [51, 52, 53] to compare item ranking approaches regarding 16 both measures (a). We also performed a one-tailed ANOVA 17 test [54] to assess differences between the three dataset sizes 18 (b). Input for the test was the relative differences of task 19

20 completion times for the two item ranking approaches.

Part 2: Assessment of User Behavior (RQ₃, RQ₄). To assess 21 22 participant behavior, two authors coded the study observation notes and extracted behavioral variables [55, 56, 57]. Coding 23 conflicts were resolved by a third author not involved in the 24 coding process. Finally, behavioral variables were reviewed by 25 an external researcher not involved in their creation. Overall, 26 we distinguished between general behaviors observed in both 27 approaches, behaviors observed only in the item-based ranking 28 approach, and unique behaviors of the attribute-based approach. 29 Re-iterating over the study observation notes, one author fur-30 ther assigned participants a score between 1 (not present) and 31 5 (very pronounced) for each behavioral variable derived. Ul-32 timately, we used the observational data, participant knowl-33 edge assessments, and behavioral variables with their manifes-34 35 tations for the identification of personas (RQ_4) . For the analysis of interactions between the two approaches (item-based and 36 attribute-based), behavioral variables, and personas, we created 37 two heatmaps shown in Figure 7. 38

8. Results of the User Study

8.1. Part 1: Performance Analysis (RQ1, RQ2)

Figure 6 (A) shows task completion time with RankASco and with the general purpose tool (RQ_1). Clearly, using the general purpose tool resulted in greater variability and outliers for participants. Conversely, independent of the size of the dataset, there is less variability regarding completion times when using RankASco. However, there was no statistically significant difference in the average task completion time (t(23) = -0.588, p = 0.563).

Figure 6 (B) includes box plots of the relative completion time difference for RankASco versus the general purpose tool, across three dataset sizes (RQ_2).

Values > 0 indicate that a participant was faster when using RankASco versus the general purpose tool and vice versa for values < 0.

Looking at the visualization depicted in Figure 6 (B), three findings stand out. First, participants were on average faster using the general purpose tool when the dataset contained 500 items (median > 0). Second, participants were faster using RankASco when using the 100 items set (median < 0). The third finding is that for 300 items dataset the task completion time was most diverse for the general purpose tool. In summary, our assumption that RankASco as an attribute-based approach performs better for larger datasets was not observed. We believe that different types of user behaviors had a stronger effect on the task completion time than the dataset size.

Figure 6 (C) reveals that it took participants longer using a 66 smaller item set (RQ_2) . One explanation of this finding could 67 be that for only 100 items, several participants did take the time 68 to traverse and interpret the entire item collection, in contrast 69 to larger item sets. Another possible explanation for this un-70 expected finding is that randomly selecting a subset from the 71 500-item dataset reduced the number of suitable apartments 72 substantially. As a result, few to no items remained after im-73 plementing all the preferences set by the Fischer family in Ex-74 cel. Confronting this problem could have reinforced the expres-75 sion of user behavior, as discussed in the next section. Some 76 participants adopted a very pragmatic approach to the issue: 77 "After realizing that there is no optimal solution that can be 78 found with the filters, I only considered the two most impor-79 tant preferences and disregarded all the others." - P12. Other 80 participants went over each item one at a time, in an effort to 81

30

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

² the two approaches and personal preferences on approach

³ usage (cf. supplemental material).

see how they could best balance these subpar results on a peritem level. We assume that this can be considered a turning 2 point that leads to the higher variance regarding participants' time in Excel, as shown in Figure 6 (C). A one-way ANOVA revealed that there was a statistically significant difference in the percent difference of completion time between at least two dataset sizes (F(2, 24) = 6.43, p < 0.01). Tukey's HSD Test for multiple comparisons showed that the mean value of the percent difference was significantly different between the 100 and 300 item datasets (p < 0.01, 95% C.I. = [0.113, 0.680]), 10 and notably different between the 300 and 500 item datasets 11 (p = 0.074, 95% C.I. = [-0.545, 0.022]). This underlines the 12 special nature of the 300-item dataset. 13

Figure 6 (D) shows the across-subject item agreement re-14 garding the top-20 ranks of participants assigned to the 100 15 item dataset (N=8) for the two ranking approaches (RQ_1) . The 16 results of the remaining participants (N=16) for the 300 and 17 500 item datasets can be found in the supplemental materi-18 als. The number on the x-axis signifies how often an item-19 was picked by participants, where 8 indicates that an item 20 was picked by all participants and one that an item was only 21 picked by one participant. Figure 6 (D) clearly shows that the 22 attribute-based approach had greater across-subject item agree-23 ment with the 100 item dataset: participants picked 100 and 24 25 44 different items for their top 20 items using the item-based and attribute-based approach, respectively, 44 items in com-26 mon between approaches. A Wilcoxon Ranked-Sum test re-27 vealed that the distribution of across-subject item agreement 28 between the ranking results of RankASco versus the general 29 purpose tool is significantly different for the 100 items dataset 30 $(U = 3.6, n_1 = 44, n_2 = 100, p < 0.01$ two-tailed). Differ-31 ences regarding the across-subject item agreement in the top 20 32 33 ranks and item count distributions were not notable for the 300 and 500 item sets (results in the supplemental material). 34

Based on our performance analysis, the selection of an appro-35 priate ranking approach should be influenced by the characteristics of the dataset, including its size and attribute composition. 37 Our findings indicate that RankASco produces superior results 38 in terms of completion time variability and across-subject item 39 agreement. The smaller completion time variability suggests 40 that RankASco may be particularly well-suited for use with a 41 heterogeneous user pool. We discovered that the difficulty of 42 the ranking task is influenced not only by the complexity of the 43 user preferences but also by the items present in the dataset at 44 hand. For multiple criteria to be considered or item sets which 45 do not match preferences nicely, users may have to start com-46 promising their preferences, which adds to the item ranking 47 challenge. In these cases, using an attribute-based approach like 48 RankASco may be preferable, as users can refrain from tedious 49 and complex item-level comparisons. 50

In summary, we found partial evidence to confirm RQ_1 and RQ_2 . However, we also learned that the behavior of individual users may have a much bigger effect on ranking effectiveness and efficiency as assumed. Informed by this finding, we next present the assessment of user behavior.

56 8.2. Part 2: Assessment of User Behavior (RQ₃, RQ₄)

57 We structure the assessment into low-level behavior of users 58 and high-level personas, as an abstraction of user behaviors.

	Behavioral variable	Description						
	Comparing Item Details	describes the degree to which participants took item details into account.						
	Neglecting Preferences	describes the number of preferences ignored by participants while creating the ranking.						
G	Query complexity	describes the degree to which filtering and sorting operations were applied.						
	Softening Preference Specifications	describes the behavior to allow flexibility for some attribute preferences.						
	Top-Rank Determination	describes the speed at which participants were willing to decide on a winning item.						
	Undo	describes the willingness of participants to refrain from actions taken and re-iterate.						
	Grouping Items	describes the strategy to create and elaborate on subsets in the items.						
	ASF Interface Exploration	describes the extent to which participants explored the design space of the given approach.						
D	ASF Creation	describes how precise participants tried to match the preferences with the ASFs.						
, r	ASF Fine-Tuning	describes how much fine-tuning participants performed when refining an ASF.						
	ASF Heterogeneity	describes the heterogeneity of ASF types the participants used to create the final ranking.						
	Neglecting Preferences	describes the degree to which participants ignored preferences given through the task.						
	ASF Resetting	describes the frequency of participants resetting an ASF to a previous state.						
	Considering Input/Output Effects	describes whether participants used the input/output distribution charts in the ASF creation process.						

Table 1. Short description of the behavioral variables (extended version, cf. supplemental material), separated by with which approach they were observed: using RankASco (R) or the general purpose tool (G).

8.2.1. Low-Level Behavioral Variables (RQ₃)

We observed many behavioral patterns in how participants 60 addressed the ranking-creation task, which we distilled into 61 behavioral variables shown in Table 1. Most participants 62 repeatedly used the task description in between the steps to 63 create the item ranking. While most participants marked 64 important statements in the narrative, some participants spent 65 a lot of time manually weighting the different attributes in 66 the task description. Only after pre-processing was completed, 67 the participants' timing started. The full description of the 68 observed behaviors and interview results can be found in the 69 supplemental material. A concise version is depicted in Table 1. 70

59

71

72

73

74

75

76

77

78

79

80

8.2.2. High-Level Personas (RQ₄)

We report on the discovery of three personas, based on observations of participants in the study. These personas are the result of coding the observational data and studying the behavioral variables [32] and their manifestations, which are depicted in Table 1. We then summarized similar manifestations of the behavioral variables, which resulted in three personas: (1) Peter, the perfectionist (N=10), (2) Eva, the explorer (N=7), and (3) Pippa, the pragmatist (N=7). In the following, we present each persona in detail.

Peter is a *perfectionist*. He strives to meet all standards of his 81 objectives with the utmost accuracy. He often and thoroughly 82 checks the original specifications of a task throughout the pro-83 cess. Also, he tends to conduct micromanagement along the 84 way, which is why time management can be problematic for Pe-85 ter. Using RankASco, his goal is to create the most precise and 86 detailed ASFs to represent the preferences of the Fischer family 87 as accurately as possible. Peter might not use all capabilities of 88 a tool, but rather optimizes the output with the capabilities that 89 he is aware of, thus going towards a local optimum. He may be-90 come discouraged if the interactive options offered, e.g., for the 91 ASF creation do not satisfy his need for perfection. For Peter, 92 using Excel for item ranking is not easy, as the means to ex-93



Fig. 7. Behavioral variable manifestations of participants using RankASco or the general purpose tool (Excel), grouped by the three personas

press general ranking preferences formally are missing. Also, the number of pairwise item comparisons needed with the gen-

eral purpose tool is a challenge, especially for increasing dataset
 sizes. Overall, we observed the tendency among perfectionists
 to prefer RankASco over Excel.

Eva is an explorer. Her goal is to try out and experiment with all capabilities of a given tool before concentrating on resolving the task at hand. Her exploratory nature helps Eva gain in-depth knowledge of the approach's design space, helping her assess which functionality is best suited to address her goal. As 10 a downside, Eva may lose sight of her goal. Using RankASco 11 as an example, Eva first experiments with all different types 12 of ASF-creation interfaces before creating actionable attribute 13 preferences. Also, Eva will reset, refine, or even undo interme-14 diate results during the process to fully investigate and finally 15 exploit the capabilities of the visual interface. We observed this 16 17 in both RankASco and Excel. Fine-tuning and achieving the most meaningful ASF is not her highest priority. Using Ex-18 cel, Eva heavily applies filtering and sorting functionality, and 19 20 she accepts that some actions may turn out not useful and need to be reverted. If this notion turns too much into a try-and er-21 ror manner, her approach may be time-consuming. As a result, 22 in time-critical situations, a less complex application may be 23 preferable to avoid distraction and to help explorers streamline 24 their actions. Overall, we observed that explorers can get lost in 25 the functionality provided by the ranking tool at hand, but they 26 27 can work effectively with both RankASco and Excel.

Pippa is a pragmatist. She wants to get things done effi-28 ciently. When problems arise, she is willing to compromise 29 on preferences and accept lower-quality task completion. Her 30 approach to problems is straightforward and linear, i.e., rather 31 less looking to the left and right. Pippa applies a clear and rig-32 orous prioritization of preferences, while possibly neglecting 33 preferences of lower priority. Pippa's approach hardly involves 34 resetting actions or decisions made, as she puts less emphasis 35 on fine-tuning, refinement, and reflection. In RankASco, she 36 initially selects the ASF that she believes to be the most prac-37 tical, e.g., Two-Point Linear, and tends to use this functional-38 ity repeatedly, even if some preferences of the Fischer family 39 would require more appropriate ASF types. When using Ex-40 41 cel, Pippa is among the fastest to complete the task, regardless of the dataset size. The reason is simple: Pippa does not sys-42 tematically inspect all items given, but is fine with seeing some 43 promising items at the top. This has a strong positive effect on 44

task completion time but a negative effect on task performance. In Excel, Pippa is also one of the first to discard preferences if they are contradicting, overly complicated, or fail to yield the desired outcomes. This type of complexity is what Pippa would like to avoid, due to her practical and pragmatic nature. Overall, we observed that using the general purpose tool was more suitable for Pippa, as she was less keen on considering multiple attributes in parallel for decision-making.

The user groups and associated personas show that using 53 RankASco is most appropriate for meticulous and perfection-54 ist users. For more pragmatic users, the general purpose tool 55 is more suitable, as the per-item operations were preferred over 56 multiple attribute-based actions needed. Finally, using the gen-57 eral purpose tool tends to be faster for exploratory users, while 58 using RankASco performs better in terms of confidence in the 59 rankings produced and approach usability, highlighting the dif-60 ference between speed and perceived success. 61

9. Discussion and Future Work

Personas. We have identified three personas based on the ob-63 servation of participants across approaches and dataset sizes. 64 As we derived the personas as a result from the higher-level 65 analysis of user behavior after the study, we did not have the 66 chance to systematically analyze personas during the study, 67 e.g., with respect to the usage of the eight types of ASFs, which 68 could be insightful. Looking forward, it would be interesting to 69 design and develop future ranking approaches with the aware-70 ness for personas. Interesting decisions include determining 71 if every persona needs its own design, or if future approaches 72 manage to incorporate and support all three personas. 73

Experiment: Selection of Approaches. We have decided for 74 RankASco as the representative of a multi-criteria attribute-75 based ranking tool and Excel as a representative of a well-76 known general purpose tool used in everyday-live situations. 77 Although this was well thought through and led to interesting 78 findings, one difference between these approaches is the novelty 79 of RankASco. In contrast to Excel, the learning curve of partici-80 pants for RankASco needed to include both tool familiarization 81 and task adoption. Beyond Excel and on the long run, different 82 tools with different support for item-based and attribute-based 83 interactions may exist, which would be worth studying.

Experiment: Study Design. We designed the experiment in a 85 way that every participant was asked to use both approaches 86 (randomized) for one pre-determined dataset size (100, 300, or 87 500 items), as the comparison between the general purpose tool 88 and RankASco was key. As an alternative, the randomized as-89 signment of users to dataset sizes 100, 300, and 500 items may 90 have revealed stronger results on the assessment and usefulness 91 of the two approaches with respect to dataset size. The assump-92 tion that using RankASco would scale better for large datasets 93 was not found, possibly due to other influencing aspects that 94 require a clearer characterization, such as the pragmatism per-95 sona or the sampling method for 100 items. Pertaining to the 96 assumption, future work includes determining the break-even 97 point where stringently attribute-based approaches outperform 98 other ranking approaches, which are less agnostic to the item 99 count. 100

46

47

48

49

50

51

52

Preprint Submitted for review / Computers & Graphics (2023)

References

No Quantitative Assessment of Accuracy. The accuracy of users when working with different approaches and dataset sizes 2 was difficult to assess quantitatively. The reason for this is the lack of ground truth information for the preference-based item ranking case, which does not allow for a quantitative perfor-5 mance evaluation in that regard. The assessment of accuracy, relevance, precision, or similar measures known in machine learning, information retrieval, and similar, is a subject of future 8 work. We identify the lack of clearly defined and formalized 9 ground truth scenarios for ranking creation and we are working 10 on methodologies to address this. 11

When am I finished?. The subjective ranking-creation task 12 based on preferences of the Fischer family is one out of many 13 possible ranking creation goals. We have observed an in-14 teresting pattern across participants. This multi-truth situa-15 tion is difficult to validate and the procedural perspective re-16 vealed challenges for many participants: when is a ranking-17 creation task finished? In general, we believe that the class of 18 preference-based creation/modeling/learning tasks may benefit 19 from process-oriented methodologies that guide designers but 20 also users through the process. 21

Task Complexity. We have assumed that the task complexity 22 would increase with the number of items involved. However, 23 during the study, we discovered that the fit of items to the rank-24 ing goal can be confounding with respect to task complexity. 25 For 100 items only, participants discovered only very few items 26 that matched the users' preferences for the ranking task. The 27 result was unexpected: users took longer to decide for the set of 28 20 (weak) items to rank on top. A recommendation would be 29 to design the items for small sample sizes in a way that the task 30 complexity does not increase due to unfortunate value distribu-31

tions in items. 32

10. Conclusion 33

We presented RankASco, a visual analytics approach for 34 the human-centered creation of item rankings. RankASco en-35 ables users to interactively express and formalize preferences 36 37 on attributes, leading to a weighted ranking of items based on multiple scores; one per attribute. RankASco is the result 38 of a two-year research project with multiple design, valida-39 tion, and reflection iterations. It builds upon a conceptual [8] 40 and a technical [9] workshop paper contribution. We compare 41 RankASco to a general purpose tool in a user study with 24 42 participants, where users were tasked to create item rankings 43 for an apartment-hunting scenario. The study involved six con-44 ditions consisting of the two different ranking approaches and 45 three different dataset sizes. During the study, we also observed 46 12 variables of user behavior and studied these behaviors with 47 respect to the two approaches (RankASco and Excel). From our 48 observations, we derived three personas as well as guidelines 49 50 on the applicability of approaches for these personas. Future work includes the expansion of the empirical work to a larger 51 participant group and to the study of more conditions, such as 52 additional ranking approaches. 53

[1] Dimara, E. Bezerianos, A. Dragicevic, P. Conceptual and method-

ological issues in evaluating multidimensional visualizations for decision support. IEEE Transactions on Visualization and Computer Graphics (TVCG) 2018;24(1):749-759. doi:10.1109/TVCG.2017.2745138. Bostandjiev, S, O'Donovan, J, Höllerer, T. Tasteweights: A visual interactive hybrid recommender system. In: ACM Conference on Rec-

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100 101

102

103

104

105

106

107

108

109

111

112

119

120

- ommender Systems. ACM; 2012, p. 35-42. doi:10.1145/2365952. 2365964. [3] Wall, E, Das, S, Chawla, R, Kalidindi, B, Brown, ET, En-
- A. Podium: Ranking data using mixed-initiative visual analytdert, ics. IEEE Transactions on Visualization and Computer Graphics (TVCG) 2018;24(1):288-297. doi:10.1109/TVCG.2017.2745078.
- Kuhlman, C, VanValkenburg, M, Doherty, D, Nurbekova, M, Deva, G, Phyo, Z, et al. Preference-driven interactive ranking system for personalized decision support. In: ACM International Conference on Information and Knowledge Management. 2018, p. 1931-1934. doi:10.1145/ 3269206.3269227.
- [5] Pereira, MM, Paulovich, FV. Rankviz: A visualization framework to assist interpretation of learning to rank algorithms. Computer Graphics Forum (CGF) 2020;93:25-38. doi:10.1016/j.cag.2020.09.017
- [6] Gratzl, S, Lex, A, Gehlenborg, N, Pfister, H, Streit, M. Lineup: Visual analysis of multi-attribute rankings. IEEE Transactions on Visualization and Computer Graphics (TVCG) 2013;19(12):2277-2286. doi:10.1109/ TVCG.2013.173
- di Sciascio, C, Sabol, V, Veas, E. urank: Exploring document recom-[7] mendations through an interactive user-driven approach. 2015,doi:10. 3140/RG.2.1.5105.0321.
- Schmid, J, Bernard, J. A Taxonomy of Attribute Scoring Functions. In: [8] EuroVis Workshop on Visual Analytics (EuroVA). Eurographics. ISBN 978-3-03868-150-2; 2021, p. 31–35. doi:10.2312/eurova.20211095.
- [9] Schmid, J, Cibulski, L, Hazwani, IA, Bernard, J. RankASco: A Visual Analytics Approach to Leverage Attribute-Based User Preferences for Item Rankings. In: EuroVis Workshop on Visual Analytics (EuroVA). The Eurographics Association; 2022, doi:10.2312/eurova.20221072. Nichols, D. Implicit rating and filtering. ERCIM; 1998,
- Jawaheer, G, Szomszor, M, Kostkova, P. Comparison of implicit and explicit feedback from an online music recommendation service. In: workshop on information heterogeneity and fusion in recommender systems. 2010, p. 47-51.
- Anand, SS, Mobasher, B. Intelligent techniques for web personalization. [12] In: Intelligent Techniques for Web Personalization. Springer; 2005, p. 1-36.
- [13] Ricci, F, Rokach, L, Shapira, B. Introduction to recommender systems handbook. In: Recommender systems handbook. Springer; 2011, p. 1-35.
- Lu, J, Wu, D, Mao, M, Wang, W, Zhang, G. Recommender system ap-[14] plication developments: a survey. Decision Support Systems 2015;74:12-
- [15] Liu, TY, et al. Learning to rank for information retrieval. Foundations and Trends in Information Retrieval 2009;3(3):225-331.
- [16] Bidoki, AMZ, Yazdani, N. Distancerank: An intelligent ranking algorithm for web pages. Information processing & management 2008;44(2):877-892.
- Zehlike, M, Bonchi, F, Castillo, C, Hajian, S, Megahed, M, Baeza-[17] Yates, R. Fa* ir: A fair top-k ranking algorithm. In: ACM on Conference on Information and Knowledge Management. 2017, p. 1569-1578
- [18] Shneiderman, B. Human-centered artificial intelligence: Reliable, safe 110 & trustworthy. Int J Hum Comput Interact 2020;36(6):495-504. doi:10. 1080/10447318.2020.1741118.
- [19] Wall, E, Das, S, Chawla, R, Kalidindi, B, Brown, ET, En-113 dert, A. Podium: Ranking data using mixed-initiative visual analyt-114 ics. IEEE Transactions on Visualization and Computer Graphics (TVCG) 115 2018;24(1):288-297. doi:10.1109/TVCG.2017.2745078. 116
- [20] Edwards, W, Barron, FH. Smarts and smarter: Improved simple meth-117 ods for multiattribute utility measurement. Organizational behavior and 118 human decision processes 1994;60(3):306-325
- [21] Tervonen, T, Figueira, JR. A survey on stochastic multicriteria acceptability analysis methods. Journal of Multi-Criteria Decision Analysis 121 2008;15(1-2):1-14.
- Cibulski, L, Mitterhofer, H, May, T, Kohlhammer, J. Paved: Pareto [22] 123 front visualization for engineering design. In: Computer Graphics Forum; 124 vol. 39. Wiley Online Library; 2020, p. 405-416. 125
- [23] Carenini, G, Loyd, J. Valuecharts: Analyzing linear models expressing 126 preferences and evaluations. In: Working Conference on Advanced Visual 127

Preprint Submitted for review/Computers & Graphics (2023)

Interfaces. ACM. ISBN 1581138679; 2004, p. 150-157. doi:10.1145/ 989863.989885.

- [24] Yuan, X, Nguyen, MX, Chen, B, Porter, DH. Hdr volvis: High dynamic range volume visualization. IEEE transactions on Visualization and Computer Graphics 2006;12(4):433-445.
- [25] Loepp, B, Herrmanny, K, Ziegler, J. Blended recommending: Integrating interactive information filtering and algorithmic recommender techniques. In: ACM Conference on Human Factors in Computing Systems. 2015, p. 975-984. doi:10.1145/2702123.2702496.

6

11

14

15

16

17

22

- [26] Pajer, S, Streit, M, Torsney-Weir, T, Spechtenhauser, F, Möller, T, Piringer, H. Weightlifter: Visual weight space exploration for 10 Weightlifter: Visual weight space exploration for multi-criteria decision making. IEEE Transactions on Visualization and Computer Graphics (TVCG) 2017;23(1):611–620. doi:10.1109/TVCG. 12 13 2016.2598589
 - [27] Isenberg, T, Isenberg, P, Chen, J, Sedlmair, M, Möller, T. A systematic review on the practice of evaluating visualization. IEEE Transactions on Visualization and Computer Graphics (TVCG) 2013;19(12):2818–2827.
- 18 [28] Carpendale, S. Evaluating information visualizations. In: Information Visualization - Human-Centered Issues and Perspectives; vol. 4950. 19 Springer; 2008, p. 19-45. doi:10.1007/978-3-540-70956-5_2. 20
- [29] Pruitt, J, Adlin, T. The persona lifecycle: keeping people in mind 21 throughout product design. Elsevier; 2010.
- 23 [30] Guo, FY, Shamdasani, S, Randall, B. Creating effective personas for product design: insights from a case study. In: Internationalization, De-24 25 sign and Global Development. Springer; 2011, p. 37-46.
- 26 [31] Revella, A. Buyer personas: how to gain insight into your customer's expectations, align your marketing strategies, and win more business. John 27 Wiley & Sons; 2015. 28
- [32] Cooper, A, Reimann, R, Cronin, D, Noessel, C. About face: the 29 30 essentials of interaction design. John Wiley & Sons; 2014.
- Johansson Fernstad, S, Jern, M, Johansson, J. Interactive quantification 31 [33] of categorical variables in mixed data sets. ISBN 978-0-7695-3268-4; 32 33 2008, p. 3-10. doi:10.1109/IV.2008.33.
- Wang, L, Sun, G, Wang, Y, Ma, J, Zhao, X, Liang, R. Afexplorer: 34 [34] Visual analysis and interactive selection of audio features. Visual Infor-35 matics 2022;6(1):47-55. doi:https://doi.org/10.1016/j.visinf. 36 37 2022.02.003.
- [35] Kuhlman, C, VanValkenburg, M, Doherty, D, Nurbekova, M, Deva, 38 39 G, Phyo, Z, et al. Preference-driven interactive ranking system for per-40 sonalized decision support. ISBN 9781450360142; 2018, doi:10.1145/ 41 3269206.3269227.
- [36] Cheng, A, Yin, Y, Yan, Z, Liu, Y, Zhou, Z. Visual analytics of multiple 42 network ranking based on structural similarity. In: 2022 IEEE 15th Pacific 43 44 Visualization Symposium (PacificVis). 2022, p. 196-200. doi:10.1109/ PacificVis53943.2022.00032. 45
- 46 [37] Carenini, G, Loyd, J. Valuecharts: analyzing linear models expressing preferences and evaluations. In: working conference on Advanced visual 47 48 interfaces. 2004, p. 150-157.
- 49 [38] McHugh, ML. The chi-square test of independence. Biochemia medica 2013;23(2):143-149. 50
- [39] Patten, ML. Understanding research methods: An overview of the essen-51 tials. Routledge; 2017. 52
- [40] Wong, BD. Bézierkurven: gezeichnet und gerechnet: Ein elementarer 53 Zugang und Anwendungen. Orell Füssli; 2003. ISBN 978-3280040218. 54
- 55 [41] Wilson, ML. Search user interface design. Synthesis lectures on infor-56 mation concepts, retrieval, and services 2011;3(3):1-143.
- [42] Davis, L. Designing a search user interface for a digital library. Jour-57 nal of the American Society for Information Science and Technology 58 2006;57(6):788-791. 59
- [43] Begel, A. Codifier: a programmer-centric search user interface. In: work-60 shop on human-computer interaction and information retrieval. 2007, p. 61 62 23 - 24.
- 63 [44] De Pauw, J, Ruymbeek, K, Goethals, B. Modelling users with item metadata for explainable and interactive recommendation, arXiv preprint 64 arXiv:220700350 2022;. 65
- [45] Petridis, S, Daskalova, N, Mennicken, S, Way, SF, Lamere, P, Thom, 66 67 J. Tastepaths: Enabling deeper exploration and understanding of personal preferences in recommender systems. In: International Conference on 68 Intelligent User Interfaces. 2022, p. 120-133. 69
- [46] Bernard, J, Steiger, M, Mittelstädt, S, Thum, S, Keim, D, Kohlhammer, 70 71 J. A survey and task-based quality assessment of static 2d colormaps. vol. 9397. SPIE Press; 2015, p. 93970M-93970M-16. doi:10.1117/ 72 12.2079841. 73
- 74 [47] Bartelheimer, Apartment rental offers in germany. 2020. https://www.kaggle.com/datasets/corrieaar/ URL

apartment-rental-offers-in-germany.

- [48] Ahlberg, C, Shneiderman, B. Visual information seeking: tight coupling of dynamic query filters with starfield displays. In: Conference on Human Factors in Computing Systems, CHI. ACM; 1994, p. 313-URL: https://doi.org/10.1145/191666.191775. doi:10. 317. 1145/191666.191775. [49] Dimara, E, Bezerianos, A, Dragicevic, P. Narratives in crowdsourced
- evaluation of visualizations: A double-edged sword? In: CHI Conference on Human Factors in Computing Systems. 2017, p. 5475-5484.
- [50] Owen, DB. The power of student's t-test. Journal of the American Statistical Association 1965;60(309):320-333.
- Hogg, RV, Tanis, EA, Zimmerman, DL. Probability and statistical [51] inference. Pearson/Prentice Hall Upper Saddle River, NJ, USA:; 2010. [52]
- McKnight, PE, Najab, J. Mann-whitney u test. The Corsini encyclopedia of psychology 2010;:1-1.
- [53] Krishnamoorthy, K. Handbook of statistical distributions with applications. Chapman and Hall/CRC; 2006.
- [54] Kim, TK. Understanding one-way anova using conceptual figures. Korean journal of anesthesiology 2017;70(1):22–26.
- [55] Charmaz, K. Constructing grounded theory: A practical guide through 95 qualitative analysis. sage; 2006. 96 Glaser, BG, Strauss, AL. The discovery of grounded theory: Strategies 97
- for qualitative research. Routledge; 2017. 98 [57] Brehmer, M, Sedlmair, M, Ingram, S, Munzner, T. Visualizing 99 dimensionally-reduced data: Interviews with analysts and a characteri-100

zation of task sequences. 2014, p. 1-8.

77

78

79

80

81

82

83

85

86

87

88

89

90

91

92

93

94

Research Highlights:

- The presentation of RankASco, an attribute-based visual analytics approach that accepts user preferences to create rankings for large item sets. RankASco is the result of a two-year research project, with two workshop paper publications forming the baseline for this extended version. We build upon RankASco with additional visual interfaces, refined design choices, and more descriptive details.
- The validation of RankASco in a user study with 24 participants. The study evaluates RankASco in comparison to Excel as a representative of a general purpose tool. We decided to recruit non-experts with low familiarity in using programmatic solutions to solve multi-criteria ranking problems.
- The presentation of three personas, characterizing common user behaviors in rankingcreation: (1) Peter, the perfectionist, (2) Eva, the explorer, and (3) Pippa, the pragmatist.

•

Overview of	f all attribute	es		Attribute So	oring Function	Resu	lt Over	/iew	Download JSON	Download Op	tions Downle	ad Weights	Download CSV	/ back
Categorical Attribut	es			- d	1	PCA of Inp	ut distribution	ł	oaseRent	condition y	/earConstructed	d living	Space PCA of	Output distribution
noParkSpaces	condition	noRooms	regio	Control Lines Mail Air Control Lines Mail Air Control Lines Mail Air Control Lines Mail Air Control Lines Control Lines	There are 0 darm with a mining rule. Being much a mining rule. Being much are of points.	-1 -3		0 -	0.85	0.71	0.28	0	40	35 1
Missing Values: 0	Missing Values: 0	Missing Values: 0	Missing Values: 0			Rank	ID	Color	item Score 🔅	Item Score Symbol	baseRent	condition	yearConstructed	livingSpace
Add	AD3	A03	A00			1	111161515		1.182		0.507	0.522	0.778	0.402
Numerical Attribute	s			42 /2 WH	-100 100 200	2	113596488		1.280		0.925	0.522	0.687	-0.178
baseRent 🔳	yearConstructed	livingSpace	serviceCharge			3	106922601		1.078		0.489	0.522	0.687	0.243
	150		0 0 1	Attribute So	oring Function	4	115370074		1.047		0.450	0.522	0.778	0.185
	100	40		Trad abstrator	Digit Definition	5	82393787		0.792		0.078	0.522	0.999	0.185
\$32 2170	1970 2020	945 200	913 332.75			6	110307393		0.909		0.346	0.522	0.778	0.064
Missing Values: 0	Missing Values: 0	Missing Values: 0	Missing Values: 0			7	106864949		0.925		0.345	0.522	0.563	0.253
Correlation between Numerical	ADJ	Correlation between Catego	Pricel Attributes (Chi Spurce)			8	114547944		0.932		0.450	0.522	0.778	-0.101
www.compet		ithindfor-		 Score Assignment Equidietant 		8	114547944		0.932		0.450	0.522	0.778	-0.101
tviegSgace-		gardenOrdaktory repo	0.6	○ Non-Equidistant		10	115106022		0.716		-0.002	0.522	0.778	0.320
yeerConstructed basefund www.ccDarge_livergiga	er yea/Carisburged Lewiferd				shitune,we shitunewated 2 4 5 (rospolered)	1 2	3 4 5 6	7 8	9 10 Next Last					
	(*	(2)					(3)						

Clara-Maria Barth: Conceptualization, Methodology, Software, Validation, Formal analysis, Investigation, Writing, Visualization, Supervision, Project administration

Jenny Schmid: Conceptualization, Methodology, Software, Investigation, Data Curation, Writing, Visualization

Ibrahim Al-Hazwani: Conceptualization, Writing

Madhav Sachdeva: Conceptualization, Software, Data Curation, Writing, Visualization

Lena Cibulski: Conceptualization, Writing

Jürgen Bernard: Conceptualization, Methodology, Validation, Resources, Writing, Supervision, Project administration, Funding acquisition

Declaration of interests

 \boxtimes The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

□The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: