

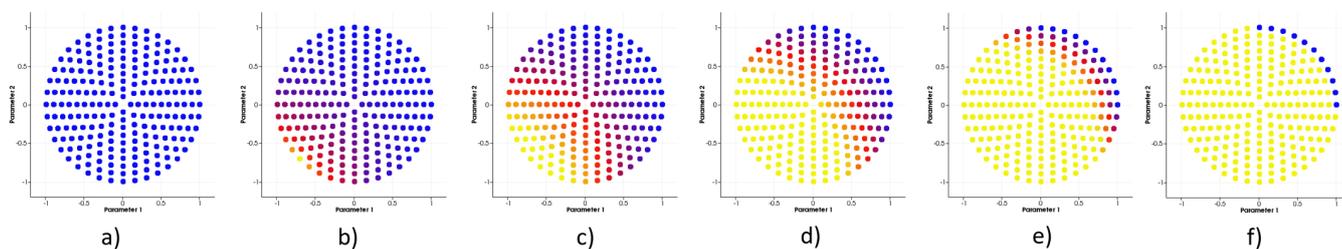
# A High-Dimensional Data Quality Metric using Pareto Optimality

Tobias Post<sup>1</sup>, Thomas Wischgoll<sup>2</sup>, Bernd Hamann<sup>3</sup> and Hans Hagen<sup>1</sup>

<sup>1</sup> Computer Graphics and HCI Group, University of Kaiserslautern, Germany

<sup>2</sup> Advanced Visual Data Analysis Group, Wright State University, U.S.A.

<sup>3</sup> Department of Computer Science, University of California (UC Davis), U.S.A.



**Figure 1:** The introduced Pareto factor is used to evaluate the quality of data points with respect to Pareto optimality. A scatterplot visualization is extended, showing the value of this new metric using a color scale ranging from yellow (0.0) over red (0.5) to blue (1.0). The effect of different Pareto exponents is shown for exponents (a) 0.0, (b) 0.25, (c) 1.0, (d) 4.0, (e) 16.0 and (f)  $\infty$ .

## Abstract

The representation of data quality within established high-dimensional data visualization techniques such as scatterplots and parallel coordinates is still an open problem. This work offers a scale-invariant measure based on Pareto optimality that is able to indicate the quality of a data point with respect to the Pareto front. In cases where datasets contain noise or parameters that cannot be expressed or evaluated mathematically, the presented measure allows to visually encode the environment of a Pareto front to enable an enhanced visual inspection.

## 1. Introduction

The representation of data quality within a high-dimensional dataset was mentioned as one of the top challenges in information visualization [LK07]. This is especially the case when the selected dataset contains noise or data dimensions that cannot easily be evaluated mathematically, such as aesthetics. Although Pareto optimality is a widely used concept to identify optimal points in a high-dimensional space, hidden dimensions justify to not only consider optimal but also nearly optimal points. So the concept of Pareto optimality has to be extended to obtain a measure on how Pareto optimal a data point is.

To address this problem, this work introduces the Pareto quotient, describing how many data points are more optimal than the considered one in the sense of Pareto optimality. It is possible to utilize the introduced Pareto quotient to evaluate data points, and therefore guide the user not only to solutions that lie on the Pareto front, but also to interesting points that are located close to the Pareto front. The visualization can be embedded into established visualization

techniques such as scatterplots or parallel coordinates in a straightforward manner, as shown in this work.

Therefore, this work contributes:

- A scale-invariant and flexible measure based on Pareto optimality, called Pareto quotient
- Visual encoding of Pareto quotients in established information visualization techniques

## 2. Related Work

Pareto optimality is a widely used concept to identify optimal high-dimensional points in an arbitrary space [EM11, FS06]. This Section summarizes visualization techniques that are based on the concept of Pareto optimality.

Different applications, such as fishery or architecture, use Pareto optimality to identify interesting data points [OKPK11, BMPM12]. Although this directly results in a set of optimal data points, these

domains are usually confronted with various dimensions not well expressible or previously unknown. Therefore, this paper extends the definition of Pareto optimality.

Ruotsalainen et al. [HRH08] use the gradient of the Pareto front to help the user navigate through this front to find interesting solutions. Although this is a suitable technique to navigate through Pareto optimal points only, hardly to express qualities like aesthetics require to extend the space of optimal solutions away from only Pareto optimal points to points that are not quite Pareto optimal, but instead have better other properties like aesthetics. This can be accomplished with the measure introduced in this work.

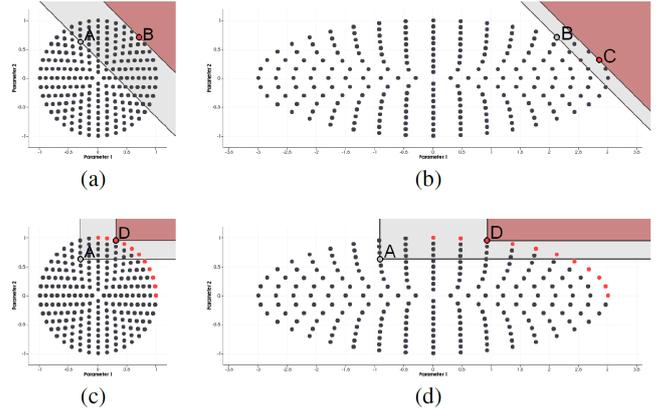
Witowski et al. [KW09] perform a study on how several known visualization techniques can be applied to visualize the Pareto front. They point out, that the combination of several tools is most promising to visualize the Pareto front in a suitable way. These techniques can be extended by the measure introduced in this work, thereby enabling users to evaluate the quality of found solutions with respect to Pareto optimality.

### 3. Methods and Results

In compensation criteria like the Kaldor-Hicks efficiency [Str] known from economics, a data point is defined to be more efficient, if the sum of all gains is greater than the sum of all losses in comparison to another data point. Here, an optimal point is a data point for that no other data point is more efficient. Fig. 2 (a) shows the optimal point *B* in red and all other points in dark gray. Point *A* is no Kaldor-Hicks optimum because there exists a more efficient point above the diagonal going through point *A*, and point *B* is optimal because there exists no such point for *B*. The problem with this kind of measure is, that it is not invariant to anisotropic scaling. Fig. 2 (b) shows the same set of data points, this time scaled anisotropically by a factor of three in the horizontal direction. Here, point *B* is not optimal any more because there exist more efficient points, and point *C* becomes the new optimum.

To avoid this problem, this work is based on Pareto optimality which is scale-invariant [Fle05]. A data point is Pareto optimal, if there exists no other data point that is greater in at least one dimension, while not being less in all other dimensions. Fig. 2 (c) and (d) show the same data points with the same scalings as before, this time evaluated with Pareto optimality. As can be seen, the red Pareto optimal points remain optimal after anisotropic scaling.

A remaining problem of Pareto optimality is, that for non optimal data points no information is provided about how far away from optimality these points are. Therefore, this work introduces a novel measure by extending Pareto optimality, thereby preserving the important invariance to anisotropic scaling. The question arises, what near or far from an optimum means in a scale-invariant way. The quotient of Pareto optimal points to all data points as described in equation (1) and (2) is used, where  $<_p$  and  $>_p$  means less and more efficient in the sense of Pareto optimality, respectively. Here,  $\vec{x}_i$  is point number *i* out of *n* different data points and  $p(\cdot)$  is the new measure called *Pareto quotient*. This measure preserves scale-invariance, what is needed if a scaling or weighting of individual dimensions in a multivariate optimization is not known or should not be determined.



**Figure 2:** A set of data points with different scalings. (b) and (d) are anisotropically scaled by a factor of three. (a) and (b) show the optimality towards a compensation criterion and (c) and (d) show Pareto optimality (red points are optimal). From (a) to (b), the optimum changes from point *B* to point *C*, meaning that this measure is dependent on the (anisotropic) scale. In contrast to that, the Pareto optimality in (c) and (d) is scale-invariant with multiple points being optimal.

Normalizing this measure to range from zero to one for all data points yields what will be called the *normalized Pareto quotient* and is useful for evaluation, weighting or visual analysis. So far, the user cannot choose how far all data points are considered or how far only Pareto optimal points are of interest. To compensate this, the normalized Pareto quotient is raised to some user defined power, the *Pareto exponent*, thereby creating what will be called the *Pareto factor*. Fig. 1 (a) - (f) show the effect of different Pareto exponents, where higher Pareto exponents focus on Pareto optimal points only, while lower Pareto exponents preserve an overview over all data points.

$$p(\vec{x}_i) = 1 - \frac{|\{j \mid j \in \{1, \dots, n\} \setminus \{i\} \wedge \vec{x}_i <_p \vec{x}_j\}|}{n-1} \quad (1)$$

$$= \frac{|\{j \mid j \in \{1, \dots, n\} \setminus \{i\} \wedge \vec{x}_i >_p \vec{x}_j\}|}{n-1} \quad (2)$$

Fig. 1 shows how the Pareto factor can be applied to established information visualization techniques like scatterplots. A color range from yellow (Pareto factor of 0.0) over red (Pareto factor of 0.5) to blue (Pareto factor of 1.0) is used. Users can interactively manipulate the Pareto exponent to focus more or less on the Pareto optimal points only. The applicability of the presented measure is not limited to scatterplots only, but can also be applied to scatterplot matrices, parallel coordinate plots, or star plots, for example. Within these techniques, the presented Pareto factor allows to evaluate the optimality of the data point.

### 4. Conclusion

This work introduced a novel measure for the quality of data points in a high-dimensional space based on Pareto optimality. It was shown that the introduced measure is scale-invariant and enables

the evaluation of the efficiency of data points with respect to Pareto optimality. Based on this measure, it was possible to visually extend established information visualization techniques. This helps to not only consider Pareto optimal data points in the analysis that might not be the desired solution for problems with noise or hard to evaluate parameters like aesthetics, but also consider nearly Pareto optimal points and find a suitable overall tradeoff solution.

## References

- [BMPM12] BOOSHEHRAN M., MÖLLER T., PETERMAN R. M., MUNZNER T.: Vismon: Facilitating analysis of trade-offs, uncertainty, and sensitivity in fisheries management decision making. URL: <http://eprints.cs.univie.ac.at/4177/>. 1
- [EM11] ESKELINEN P., MIETTINEN K.: Trade-off analysis approach for interactive nonlinear multiobjective optimization. *OR Spectrum* (2011). 1
- [Fle05] FLEISCHER M.: Scale invariant pareto optimality: A meta-formalism for characterizing and modeling cooperativity in evolutionary systems. In *Proceedings of the 7th Annual Conference on Genetic and Evolutionary Computation* (2005), GECCO '05, ACM, pp. 233–240. 2
- [FS06] FELDMAN A. M., SERRANO R.: *Welfare Economics and Social Choice Theory*, 2. ed. Springere, 2006. 1
- [HRH08] HENRI RUOTSALAINEN E. M., HÄMÄLÄINEN J.: Navigation on a pareto-optimal front utilizing gradient in formation in interactive multiobjective optimization. *International Conference on Engineering Optimization* (2008). 2
- [KW09] KATHARINA WITOWSKI MARTIN LIEBSCHER T. G.: Decision making in multi-objective optimization for industrial applications - data mining and visualization of pareto data. *European LS-DYNA Conference* (2009). 2
- [LK07] LARAMEE R. S., KOSARA R.: *Challenges and Unsolved Problems*. Springer Berlin Heidelberg, 2007, pp. 231–254. 1
- [OKPK11] OH S., KIM Y., PARK C., KIM I.: Process-driven bim-based optimal design using integration of energyplus, genetic algorithm, and pareto optimality. *Proceedings of the IBPSA building simulation 2011 conference, Sydney, Australia* (2011), 894–901. 1
- [Str] STRINGHAM E.: Kaldor-hicks efficiency and the problem of central planning. the quarterly. *Journal of Austrian Economics*, 41–50. 2