



**TO: CORRESPONDING AUTHOR**

**AUTHOR QUERIES - TO BE ANSWERED BY THE  
AUTHOR**

**Dear Author**

**Please address all the numbered queries on this page which are clearly identified on the proof for your convenience.**

**Thank you for your cooperation**

Q1	Should RELIEFF be set as ReliefF as in other sources? See throughout article.	
Q2	Please supply date of meeting, place of publication, name of publisher and page numbers.	
Q3	Please supply date and place of symposium, place of publication and name of publisher.	
Q4	Please clarify publication details.	
Q5	Please supply place of publication.	
Q6	Please supply place of publication.	
Q7	Please supply date of conference, place of publication, name of publisher and page numbers.	
Q8	Please supply date and place of conference.	
Q9	Please supply date and place of conference, place of publication and name of publisher.	
Q10	Please supply date and place of conference, place of publication, name of publisher and page numbers.	
Q11	Please supply date and place of conference, place of publication, name of publisher and page numbers.	
Q12	Please supply place of publication and name of publisher.	
Q13	Please supply date and place of conference, place of publication, name of publisher and page numbers.	
Q14	Please supply date of conference, place of publication and name of publisher.	
Q15	Please supply date of conference, place of publication, name of publisher and page numbers.	

Q16	Please supply date of conference, place of publication, name of publisher and page numbers.	
Q17	Please supply page numbers.	
Q18	Please supply place of publication.	
Q19	Please provide Table 4 citation.	

Production Editorial Department, Taylor & Francis  
4 Park Square, Milton Park, Abingdon OX14 4RN

Telephone: +44 (0) 2070 176000

Facsimile: +44 (0) 2070 176336

# A RELIEFF attribute weighting and X-means clustering methodology for top-down product family optimization<sup>1</sup>

Q1

Conrad S. Tucker, Harrison M. Kim\*, Douglas E. Barker and Yuanhui Zhang

*University of Illinois, Urbana-Champaign, Illinois, 61801*

*(Received)*

This article proposes a top-down product family design methodology that enables product design engineers to identify the optimal number of product architectures directly from the customer preference data set by employing data mining attribute weighting and clustering techniques. The methodology also presents an efficient component sharing strategy to aid in product family commonality decisions. Two key data mining models are presented in this work to help guide the product design process: (1) RELIEFF attribute weighting technique that identifies and ranks product attributes, and (2) the X-means clustering approach that autonomously identifies the optimal number of candidate products. Product family commonality decisions are guided by once again employing the X-means clustering technique, this time to identify the components across product families that are most similar. A family of prototype aerodynamic air particle separators is used to evaluate the efficiency and validity of the proposed product family design methodology.

**Keywords:** data mining; X-means clustering; RELIEFF; bi-level quasi-separable problem; product architecture; aerodynamic particle separator

## Nomenclature

AF	Air flow area.
$f_k$	Local product design objective function(s), a function of local design variables: $f_k(\mathbf{x}_k)$ .
$\mathbf{R}^{Eng}$	Engineering design response (feasible/infeasible).
$\mathbf{T}^{Cj}$	Vector of product attributes represented by the cluster centroid in the data mining model.
$\mathbf{y}_{s,k}$	Linking variable at the engineering subsystem level cascaded up to system level.
$\varepsilon_y$	Deviation tolerance between linking variables.
$\xi$	Particle separation efficiency.
$\mathbf{w}'$	The vector of newly transformed RELIEFF weights for target vector $\mathbf{T}^{Cj}$ .
$\ \cdot\ _2^2$	Squared L-2 norm notation measuring the deviation between targets and responses.
$k$	$k$ th candidate product architecture determined by the results of the X-means clustering.
$K$	The total number of cluster centroids $C_j$ that exist for the X-means clustering solution.
$Cost_k$	Total product cost represented as the summation of individual component costs.

\*Corresponding author. Email: hmkim@illinois.edu

## 1. Introduction

In complex engineering systems that require a wide range of operating conditions, engineers are left with the challenge of designing product portfolios that meet customer preferences. The product family paradigm has been proposed to address the challenges of designing products for mass customization or for highly diversified customer functionality requirements. The term ‘product family’ is frequently defined in literature as a group of related products that share an underlying product design architecture (Messac *et al.* 2002, Alizon *et al.* 2007, Tucker and Kim 2007). The product family paradigm enables companies to standardize certain aspects of a product and at the same time provide product diversity to customers through product variants. Product cost savings may be realized as a result of product standardization due to ‘economies of scale’ (*e.g.* cost savings due to a standard manufacturing line for all products, rather than a specialized manufacturing line for each product). However, greater product standardization may also lead to lower product diversity in the market space and diminished product performance (*e.g.* limited customizable features for customers such as product colour, reliability, size, etc.) Therefore, in the product family approach, the level of product standardization versus product variety presents a trade-off scenario as product performance and appeal (from the customer’s perspective) may diminish in an attempt to increase product standardization (Messac *et al.* 2002).

The product family design problem has been segmented into two well established domains: the ‘Bottom-Up’ approach and the ‘Top-Down’ approach (Alizon *et al.* 2007). In the ‘Bottom-Up’ approach, companies are more interested in making significant improvements to an existing product portfolio by combining products within the existing product family into a new product architecture. An assumption in the ‘Bottom-Up’ approach is that the newly redesigned product family will be able to satisfy customer needs through minimal additional technology investments. On the other hand, in the ‘Top-Down’ approach to product family design, the next generation of products is not based on an existing product family, but instead emerges from a market driven need. This need arises from the evolution of customer preferences far beyond what the current product portfolio can satisfy (Alizon *et al.* 2007). In this work, a ‘Top-Down’ product family methodology is proposed that analyses large customer preference data sets and identifies candidate product architectures that will be used in the product family design. This product design architecture can represent a group of design components that perform a series of functional processes. Products sharing similar product architecture can satisfy a broader range of customer requirements simply by possessing functionality capabilities that vary beyond the underlying architecture. The sharing of components also has the potential to reduce the time and costs associated with manufacturing diverse products. The challenge of developing a product architecture to be used in a product family presents an interesting design problem as the customer pool and functionality demands increase.

As data storage and information retrieval capabilities become more widely available, there is an emerging trend for companies to acquire and store customer preference data. For example, the physical characteristics (vehicle horsepower, number of doors, colour, etc.) of an automobile purchased by a customer visiting a dealership, along with the customer’s demographic information (age, gender, household income, etc.) can enable auto manufacturers to determine emerging trends in the automotive industry and design next generation products accordingly. A great challenge in storing such data for product design purposes, however, is the non-homogeneity of customers, along with their individual preferences. Therefore, as the size of this non-homogeneous data increases, so does the complexity of identifying natural patterns within the data set. The ability to determine suitable product architectures for a particular group of customers becomes a challenge as enterprise decision makers and engineers attempt to extract meaningful patterns within the data set to aid in the product design and development process. Data mining in the context of product development is an emerging area of research that has the potential to significantly impact engineering design and manufacturing efforts (Kusiak 2006, Tucker and Kim 2007). By

101 identifying patterns within the large data set of customer preferences, engineers can incorporate  
102 this knowledge in the product family design process.

103 The product family design scenario that this article focuses on is described as follows. An  
104 enterprise is launching a portfolio of products that potentially share a subset of components.  
105 However, there are a few issues that must be resolved *prior to* the design process. First, the data  
106 set used in this product portfolio design scenario comprises large-scale non-homogeneous data  
107 which indicates that the product (the aerodynamic particle separator) undergoes a wide range of  
108 operating/environmental conditions. Second, the enterprise does not have a prior knowledge as  
109 to how many product variants should be introduced in the market, although customer preference  
110 data is available from survey or market research. Third, the enterprise does not know which  
111 components should be shared in case more than one product variant is introduced, although it has  
112 the flexibility to accommodate the component sharing decisions.

113 This article presents a product family design methodology that is driven by data mining capabilities,  
114 which resolves the product family design challenges presented above. Often in the engineering  
115 design process, quantifying attribute importance from a customer's perspective is challenging.  
116 Possessing a mechanism that can identify which performance attributes are more dominantly  
117 represented in the preference data set would help the engineering design teams focus resources  
118 in a more efficient manner. For this, the RELIEFF (Kira and Rendell 1992) attribute weighting  
119 algorithm is employed to identify attributes in order of importance in the data set. Then, the X-  
120 means clustering algorithm is employed to identify groups of similar operating states within the  
121 raw data set. As a result, the number of product variants that should be introduced to reflect preferences  
122 (represented in the data) can be identified. Finally, the X-means clustering is applied again  
123 to the detailed designs of the initial product variants to identify which components may be shared  
124 among them. These sharing design decisions are implemented in a multi-disciplinary design optimization  
125 framework where an individual product variant is modelled as an individual subsystem.

126 The rest of the article is organized as follows. Section 2 provides research background followed  
127 by the proposed methodology in Section 3; the methodology is demonstrated in a case study in  
128 Section 4 followed by results and discussion in Section 5 and conclusion in Section 6.

## 131 2. Research background

132 A selective literature review on research areas pertaining to the concepts and techniques proposed  
133 in this work will be presented. These research areas were reviewed in a selective manner based  
134 on their relevance to data mining in product family design/product portfolio optimization.

### 137 2.1. Data mining in product design

138 There have been several researchers in the product design community that have incorporated data  
139 mining techniques in the product design process. For example, Agard and Kusiak (2004) utilize  
140 data mining clustering techniques to address the customer segmentation problem by determining  
141 a target market in a new product development process. Association rule mining is then used to  
142 discover attribute patterns in the segmented data (Agard and Kusiak 2004). Later works by Kusiak  
143 illustrate the benefits of data mining in a wide array of diversified industries such as biotechnology,  
144 energy, pharmaceutical, etc. (Kusiak 2006).

145 Tucker and Kim have incorporated data mining techniques in the product portfolio formulation  
146 process for extremely volatile markets (Tucker and Kim 2007, 2008). In such industries, product  
147 life cycles are short lived. Therefore, being able to correctly predict a customer's product preferences  
148 is paramount to increasing a product portfolio's chances of market success. Tucker and  
149  
150

151 Kim (2007) approached this design problem by systematically linking a customer's preferences,  
 152 acquired through predictive data mining techniques, directly with engineering detailed design  
 153 through multilevel optimization techniques (Kim *et al.* 2002, Kim *et al.* 2003).

154 Data mining techniques are also employed by Moon *et al.* (2006) in representing the functional  
 155 requirements of customers. The proposed methodology uses fuzzy clustering techniques to deter-  
 156 mine the module composition of a product architecture (Moon *et al.* 2006). The work assumes that  
 157 a product is an amalgam of module-based components with prior knowledge of the functionality  
 158 capabilities of each module.

159 The primary contribution of this work is to present a product family design methodology for  
 160 complex engineering systems that autonomously identifies the number of products to design by  
 161 extracting weighted product preference information from a customer data set. This work focuses  
 162 on design problems with large product preference data sets that can be integrated into the product  
 163 design process. Since it would be impractical and highly expensive (from a cost and logistics  
 164 standpoint) to design an independent system for each operating scenario, the proposed methodol-  
 165 ogy instead identifies the most similar operating requirements given a large data set of operating  
 166 conditions/scenarios. This in turn highlights the cost savings associated with product platform  
 167 design through the concept of component sharing. This is modelled by the shared linking variable  
 168 in the bi-level quasi-separable problem formulation that attempts to achieve an optimal design  
 169 solution for each product while concurrently satisfying specific product functionality requirements  
 170 (Tosserams *et al.* 2007). The term '*quasi-separable*' is used in this work to denote independent  
 171 sub-problems that share a common design variable/component. In this work, sub-problem simply  
 172 means a unique product design. The bi-level formulation is used in this work to co-ordinate these  
 173 sharing decisions among sub-problems. Therefore, the individual sub-problem formulation for  
 174 the bi-level quasi-separable problem is as follows:

175 Minimize

$$176 \quad f_k(\mathbf{y}_{s,k}, \mathbf{x}_k) \quad (1)$$

177 Subject to:

$$179 \quad g_k(\mathbf{y}_{s,k}, \mathbf{x}_k) \leq 0 \quad (2)$$

$$181 \quad h_k(\mathbf{y}_{s,k}, \mathbf{x}_k) = 0 \quad (3)$$

182 For the quasi-separable formulation, each  $\mathbf{x}_k$  represents the vector of local design variables unique  
 183 to each sub-problem ( $k$ ), where  $k = 1, \dots, K$  sub-problems. The vector of linking variables  $\mathbf{y}_{s,k}$   
 184 makes the sub-problems quasi-separable as each sub-problem sharing a linking variable becomes  
 185 influenced by the solution of other sub-problems sharing the same linking variable. A master  
 186 problem is used to co-ordinate the linking variable among sub-problems and is explained in more  
 187 detail in Section 3.2 of this work.  
 188

## 189 2.2. Product family optimization

191 The product family design paradigm has been investigated extensively throughout the engineer-  
 192 ing design community. Although there are a wide range of application areas, the underlying  
 193 focus of product family optimization is to design a group of related products built around a  
 194 common functional system architecture/platform. The aim is that commonality among prod-  
 195 uct variants will reduce product design and manufacturing costs while still satisfying customer  
 196 requirements. There have been many proposed methodologies and metrics for evaluating product  
 197 commonality decisions in product family optimization. For example, the degree of commonal-  
 198 ity index (DCI) proposed by Collier (1981) measures the ratio of common components existing  
 199 among products within a product family to the total number of components (Collier 1981). Later  
 200 proposed commonality strategies such as the total constant commonality index (TCCI) (Wacker

201 and Trelevan 1986), the commonality index (CI) (Martin and Ishii 1996, 1997, Khajavirad and  
202 Michalek 2007), component part commonality index (CI<sup>(C)</sup>) (Jiao and Tseng 2000), product line  
203 commonality index (PCI) (Kota *et al.* 2000), the percent commonality index (%C) (Siddique *et*  
204 *al.* 1998), the generational variety index (GVI) (Martin and Ishii 2002), the functional similarity  
205 index (FSI) (McAdams *et al.* 1999, McAdams and Wood 2002), and the comprehensive metric  
206 for commonality (CMC) (Thevenot and Simpson 2007), propose strategies to help improve prod-  
207 uct commonality decisions by either rewarding or penalizing component sharing decisions. The  
208 aforementioned commonality indices are referenced in this work to give the reader a glimpse at  
209 the myriad of approaches available to address the issue of commonality in product design and  
210 development and how the proposed approach differs from them.

211 Instead of employing traditional commonality indices such as those listed above, product com-  
212 monality decisions are investigated by employing the X-means clustering technique during the  
213 product family optimization process to identify similar components among product designs,  
214 hereby avoiding the need to exhaustively search all possible component sharing possibilities.  
215 In the aerodynamic particle separator problem that is investigated, the X-means clustering tech-  
216 nique is first used to identify similarities among unique operating requirements. These clusters  
217 will form the basis of the individual product platform. For the aerodynamic particle separator  
218 problem, commonality decisions will be based primarily on the manufacturing costs associated  
219 with each unique design. The costs savings benefits of incorporating commonality decisions in  
220 the product family design process will be presented later in this work.

### 223 3. Methodology

224  
225 Figure 1 is a flow diagram visually illustrating the sequence of the proposed product family design  
226 methodology. Figure 1 begins with the acquisition of raw customer product preference data and  
227 employs data mining attribute weighting and clustering techniques to determine the number of  
228 unique products needed for a given data set. One of the novel contributions proposed in this  
229 work, to solve the top-down product family research problem, is the ability to identify the optimal  
230 number of product architectures based solely on the data set. For products with highly diverse  
231 operating conditions, the data set itself may be highly heterogeneous making it quite difficult for  
232 engineers to determine the number of products to design in order to satisfy the market space. By  
233 employing the data mining RELIEFF attribute weighting and X-means clustering techniques to  
234 the raw data set (Figure 1), engineers can determine the initial product architectures to design.

235 Steps 2 and 3 in Figure 1 illustrate the added benefits of component sharing by clustering  
236 similar products together in an attempt to reduce product design costs. The X-means clustering  
237 technique is employed at the engineering design level to determine which products are similar  
238 enough to potentially benefit from component sharing decisions. The details of Figure 1 will now  
239 be explained in depth in the following sections.

#### 242 3.1. Data mining product preferences

243  
244 The data mining of product preferences is the stage where dominant patterns are identified within  
245 the raw data set (Fayyad *et al.* 1996). With each unique instance in the data set representing a  
246 customer's preferred operating state for the product, the number of operating states can increase  
247 rapidly, thereby making it impractical for a single design to exist for each unique state. The  
248 engineering design goal is to identify those operating states within the data set that are similar in  
249 design requirements (as determined by the data mining algorithm). Since product attributes may  
250 vary in terms of design significance, an appropriate attribute weighting technique would help guide  
the engineering design process. To accomplish these product design challenges, the RELIEFF

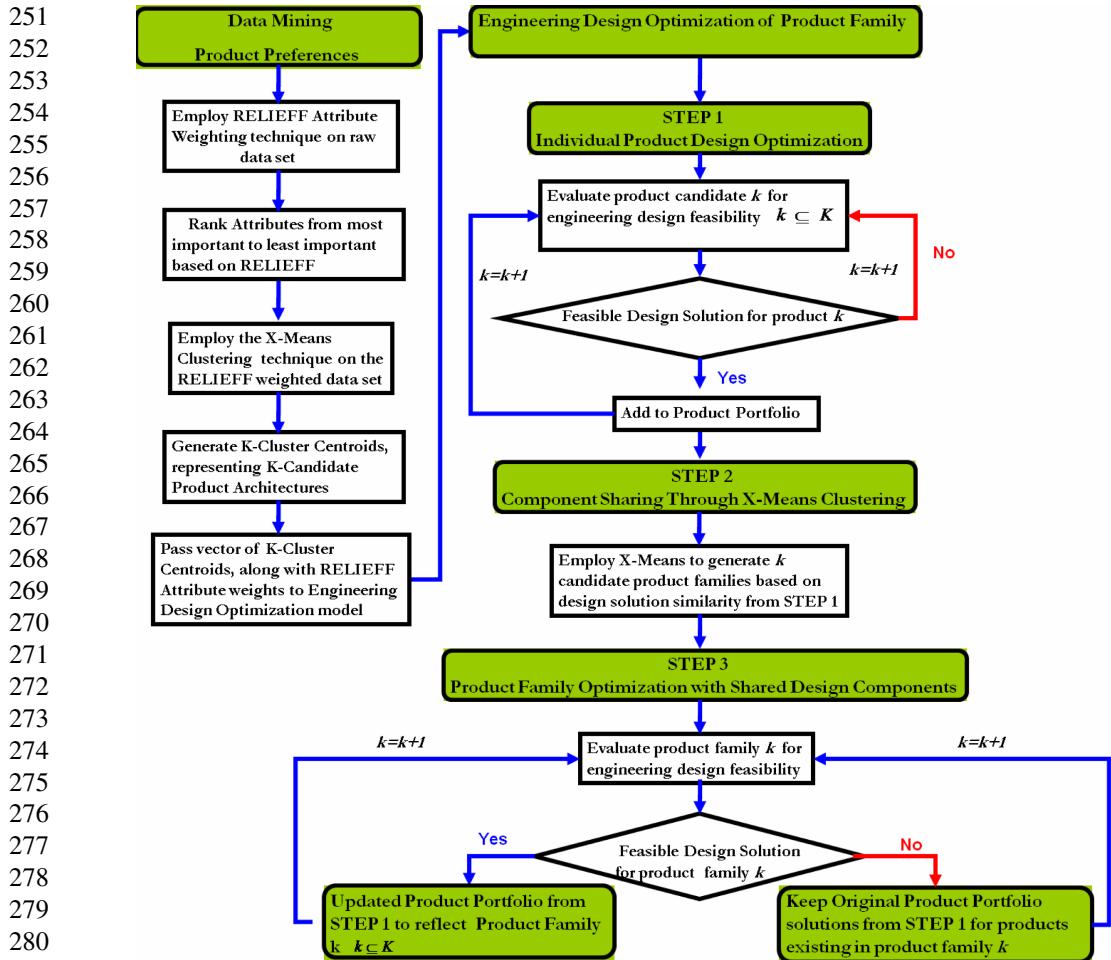


Figure 1. Flow diagram of proposed product family design methodology: From data mining product architecture identification to component sharing through X-means clustering.

attribute weighting algorithm is first employed to weight attributes in order of importance in the data set (Kira and Rendell 1992). Then the X-means clustering algorithm is used to identify groups of similar operating states within the raw data set (illustrated in the data mining flow diagram in Figure 1 and visually represented on the left in Figure 2). The weighted attributes will influence both the data clustering process as well as the engineering design model as more valued product attributes will be given more weight in the overall product family design methodology. The X-means clustering algorithm is employed again in the component sharing decisions during the product family optimization stage as similar individual product designs are grouped together by similarity of design (Steps 2 and 3 in Figure 1 and visually represented on the right in Figure 2). Below is an introduction to the RELIEFF attribute weighting technique that will later be applied to the raw data set.

### 3.1.1. RELIEFF product attribute weighting algorithm

In this work, the enhanced version of the RELIEF algorithm called RELIEFF is employed (Kononenko 1994). RELIEFF extends the original RELIEF algorithm by enabling it to efficiently

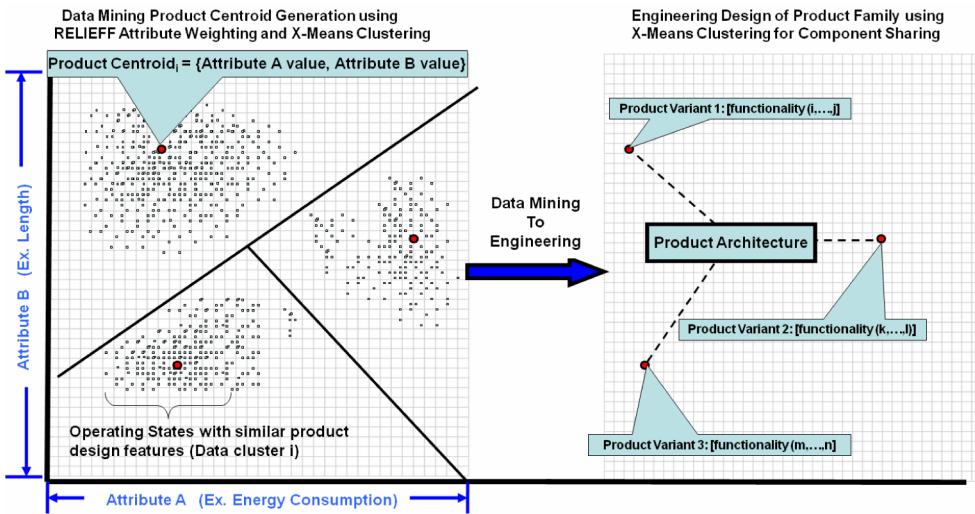


Figure 2. Visual representation of family design based on data mining RELIEFF attribute weighting and X-means clustering for product centroid generation and X-means clustering for product family component sharing optimization.

handle multi-class variables and also missing values within the data set (Kononenko 1994). A class variable can be thought of as the response or predictor variable of interest. Examples of class variables in product design data sets may include efficiency, energy consumption, price, etc.

The original RELIEF algorithm proposed by Kira and Rendell (1992) is an attribute evaluation technique that will enable product development engineers extract the importance of individual product attributes within a raw data set without explicit user provided ranking information (Kira and Rendell 1992). This can prove to be a vital time saving strategy, especially for extremely large data sets containing many attributes. Identifying the order of attribute relevance within a data set can reduce the overall computational complexity and increase the efficiency of data mining algorithms (Kira and Rendell 1992, Kononenko 1994).

Given a raw data set  $S$ ,  $m$  instances are selected to serve as the number of sampled instances where  $p$  denotes the total unique attributes within the sample set  $m$  (Kira and Rendell 1992). The overall objective of the RELIEF algorithm is to take a random sample, and using a nearest neighbour search, to identify an identical class variable, which is defined as a NEAREST HIT ( $H$ ), and also a different class variable that is nearest to the sample, defined as a NEAREST MISS ( $M$ ) (Kira and Rendell 1992). The iterative process of RELIEF estimates attribute weights  $W[A_i]$  based on their similarity to a given class, where  $A_i$ , represents a unique attribute within the data set. The general form of the algorithm can be represented as follows (Kira and Rendell 1992).

Given  $m$ , the desired number of sampled instances, and  $p$ , the number of attributes,

1. set all weights  $W[A_i] := 0.0$ ;
2. for  $j := 1$  to  $m$  do begin;
3. randomly select an instance  $X$ ;
4. find nearest hit  $H$  and nearest miss  $M$ ;
5. for  $i := 1$  to  $k$  do begin;
6.  $W[A_i] := W[A_i] - \text{diff}(A_i; X; H)/m + \text{diff}(A_i; X; M)/m$ ;
7. end;
8. end;

where the *diff* function above measures the difference between the attribute being evaluated  $A_i$  taken from the randomly selected instance  $X$ , and the value of that same attribute given the closest hit ( $H$ ) or closest miss ( $M$ ). For discrete attributes, if the value of the attribute ( $A_i$ ) of the randomly selected instance ( $X$ ) matches that of the nearest hit ( $H$ ) or nearest miss ( $M$ ), then the *diff* value is 0 (meaning values are identical), otherwise (1) meaning they are different. For continuous attributes, the actual difference is used and then normalized on a scale of [0, 1].

By its design, the RELIEFF attribute weighting technique does not constrain the attributes to non-negative values. Therefore, the weights will first be normalized based on the mini-max normalization (Han and Kamber 2006). For a given vector of attribute weights  $[w_1, w_2, \dots, w_p]$ , the weights are normalized using the following formulation:

For weight  $i = 1, \dots, p$

$$w'_i = \frac{w_i - \min_w}{\max_w - \min_w} (new\_max_w - new\_min_w) + new\_min_w \quad (4)$$

Here,

$w'_i$	The newly transformed weight ( $i$ ) of attribute ( $i$ ).
$\min_w$	The minimum value in the vector of RELIEFF attribute weights.
$\max_w$	The maximum value in the vector of RELIEFF attribute weights.
$new\_max_w$	The maximum value of the new range.
$new\_min_w$	The minimum value of the new range.

The new vector of weights  $\mathbf{w}'$  determines the level of importance for each target vector ( $\mathbf{T}^{Cj}$ ) at the engineering design level.

The data set with the newly updated weighted attributes will be used to:

- Weight clusters generated by the X-means clustering approach (discussed in the following section).
- Serve as attribute target weights for the engineering product design model.

### 3.1.2. X-means clustering

The X-means clustering algorithm in data mining is an enhancement of the k-means clustering algorithm (Pelleg and Moore 2000). Before investigating the X-means clustering algorithm and its significance in product family optimization, the k-means algorithm will be briefly described.

Given a raw data set of unique customer preferences (operating conditions), the k-means algorithm attempts to partition the original data set into  $k$  subsets of the data, where  $k$  represents the number of unique subsets or in the appropriate data mining terminology, clusters (Hartigan and Wong 1979, Jain and Dubes 1988). Each cluster contains a centroid, with data points of the cluster associated with this centroid. It is important to note that the number of clusters in the  $k$ -means algorithm is given *a priori* as a user defined input. In the context of product family design, this would be analogous to the engineering design team specifying the number of product platforms that the customer's product operating requirements must adhere to. Rather than design teams making postulations about the raw data set, a more natural process would be for the inherent patterns of the raw data set to help guide the product platform number (this is one of the contributions of the X-means data mining technique). Although there have been many enhancements to the  $k$ -means since its conception (Arora *et al.* 1998, Kanungo *et al.* 2002, Tarpey 2007), the

401 basic underlying mathematical formulation can be represented as follows:

$$402 \quad 403 \quad 404 \quad 405 \quad f = \sum_{j=1}^K \sum_{x_i \in S_j} \|x_i - c_j\|^2 \quad (5)$$

406 Here,  $S_j$  is a cluster of data points.

407 Here,  $S$  will be defined as all instances in the raw data set and, therefore,  $S_j$  would simply be a  
408 subset of this.

409  $c_j$  is the centroid of a cluster  $S_j$ .

410  $x_i$  is a data point existing within a cluster.

411  $K$  is the total number of clusters (specified a priori by the user).

412 The iterative process of the k-means algorithm begins by initially selecting the desired number  
413 of clusters ( $S_j$ ) and making an initial guess of the cluster centroid values ( $c_j$ ) (Hartigan and Wong  
414 1979). The next stage involves assigning a data point to the closest cluster centroid and centroid  
415 value (if necessary) by minimizing the error function in Equation (5) until negligible deviation  
416 occurs with each iteration.

417 The X-means clustering algorithm aims to improve on three key areas of the k-means algorithm  
418 (Pelleg and Moore 2000).

419 (1) Eliminating the need for number of clusters to be known as *a priori*.

420 (2) Improving the computational scalability.

421 (3) Enhancing the search criteria for updating cluster centroids.

422 The process by which X-means achieves these improvements is in part based on its selection  
423 criterion to determine when to add or replace a specific cluster centroid with child centroids. Child  
424 centroids originate from splitting the original solution of a k-means iteration and determining if  
425 the child clusters more accurately represent the data points once belonging to the parent centroid  
426 (Pelleg and Moore 2000). The posterior probabilities will be used to rank the models  $\Pr[M_j|D]$ ,  
427 where  $D$  represents the given data set and  $M_j$  represents each model with a given cluster size  
428  $k$ . The Bayesian information criterion (BIC) is used by X-means to rank which model is a more  
429 accurate representation of the original raw data set. Mathematically, the BIC is represented as  
430 follows (Kass and Wasserman 1995, Pelleg and Moore 2000):

$$431 \quad 432 \quad 433 \quad 434 \quad 435 \quad BIC(M_j) = l_j(D) + \frac{p_j}{2} \log R \quad (6)$$

436 Here,

437  $l_j(D)$  is the log likelihood of the data taken at the maximum likelihood point.

438  $D$  represents the given data set.

439  $p_j$  represents the number of parameters in  $M_j$ .

440  $R$  is the total number of data points of candidate centroids.

### 441 442 443 3.1.3. *Relevance of X-means to engineering product architecture design*

444 Engineering design problems involving a wide range of operating states specified by customers,  
445 can benefit from X-means clustering by identifying appropriate product functionality criterion  
446 for developing a product architecture and subsequent product family. The X-means clustering  
447 technique eliminates the need to guess the number of product architectures needed for a partic-  
448 ular customer pool by analytically generating the appropriate number of clusters (product  
449 architectures) with corresponding product functionality specifications. A user instead specifies a  
450

451 broad range for the number of clusters and X-means will identify the optimal cluster, given the  
 452 natural patterns within the data set (Pelleg and Moore 2000). This will ensure that the result-  
 453 ing product family will be a true representation of the data set for which the designs are made.  
 454 Figure 2 illustrates how the cluster centroids of the X-means data mining clustering approach  
 455 are integrated into the engineering design. The product centroids illustrated in Figure 2 rep-  
 456 resent the individual vectors of attribute value solutions that best describe *similar* groups of  
 457 customers within the raw data. Each unique product centroid will form a vector of product  
 458 preference targets used to guide the product architecture optimization process. The engineer-  
 459 ing design illustrated in Figure 2 represents the design of individual products based on the  
 460 X-means cluster centroids where each product will have unique functionality characteristics that  
 461 aim to satisfy the overall customer preference targets. Section 3.2.3 describes how product vari-  
 462 ants are then designed based on underlying product architecture under the notion of component  
 463 sharing.

### 465 3.2. Engineering design optimization of product family

#### 467 3.2.1. Step 1: Individual product design optimization

469 The results from the data mining stage provide product design engineers with several vital pieces  
 470 of information. First, the results from the X-means clustering represent the vector of product  
 471 attributes that form the product design targets ( $\mathbf{T}^{Cj}$ ) around which a product architecture is  
 472 designed (Step 1 in Figure 1). Product design targets can range anywhere from physical prod-  
 473 uct dimension targets such as length or width to product performance targets such as efficiency  
 474 or speed.

475 The second vital piece of information from the data mining stage is the relevance of each  
 476 attribute target to the customer as determined by the RELIEFF attribute weighting technique.  
 477 That is, for each attribute target vector ( $\mathbf{T}^{Cj}$ ), there will be an accompanying vector of attribute  
 478 target weights  $\mathbf{w}'$ . The engineering product architecture optimization is comprised of the detailed  
 479 engineering design model and incorporates the results from the data mining stage that help guide  
 480 the product architecture design. Here, local design variables are used to model the physical  
 481 dimensions and performance objectives of the product architecture subject to engineering design  
 482 constraints.

483 The general mathematical model for the engineering product architecture optimization is  
 484 as follows:

485 *Note:* The deviation is measured by the squared L-2 norm, which will be used throughout the engineering optimization  
 486 models presented in this work. For example:

$$488 \|\mathbf{x} - \mathbf{y}\|_2^2 = \sum_i (x_i - y_i)^2.$$

491 For the  $k$ th product architecture,

492 Minimize

$$494 F(x)_{Architecture(k)} = f_k + \mathbf{w}' \left\| \mathbf{T}^{Cj} - \mathbf{R}_k^{Eng} \right\|_2^2 \quad (7)$$

496 Subject to:

$$498 \mathbf{g}_k(\mathbf{x}_{k,}) \leq \mathbf{0}$$

$$500 \mathbf{h}_k(\mathbf{x}_{k,}) = \mathbf{0}$$

501 Here,

- 502
- 503  $f_k$  Local product design objective function (s), a function of local design variables:  $f_k(\mathbf{x}_k)$ .
- 504  $\mathbf{T}^{C_j}$  Vector of product attributes represented by the cluster centroid in the data mining  
505 model. That is, for cluster centroid  $C_j = [A_1, A_2, \dots, A_p]$ , target  $\mathbf{T}^{C_j}$  is set as  $[A_1,$   
506  $A_2, \dots, A_p]$  where  $A_1, A_2, \dots, A_p$  represent attribute values for a given centroid  $C_j$ .
- 507  $\mathbf{w}'$  The vector of newly transformed RELIEFF weights for target vector  $\mathbf{T}^{C_j}$ .
- 508  $\mathbf{R}_k^{Eng}$  Vector of engineering responses based on the formulation of the engineering design  
509 model.  $\mathbf{R}_k^{Eng}$  is a function of local design variables  $\mathbf{x}_k$  and is represented by  $\mathbf{R}_k^{Eng}$   
510  $(\mathbf{x}_k)$ .
- 511  $\mathbf{g}_k$  Inequality design constraints bounding the product architecture model.
- 512  $\mathbf{h}_k$  Equality design constraints bounding the product architecture model.
- 513  $K$  The  $k$ th candidate product architecture determined by the results of the X-means  
514 clustering.
- 515  $K$  The total number of cluster centroids  $C_j$  that exist for the X-means clustering solution.  
516

517

518 *Note:* It is important to note that although there may be  $K$  candidate product architectures to investigate, there may not  
519 always be a feasible design solution for the  $k$ th product architecture as generated product preference requirements may  
520 be too demanding, given the constraints of the engineering design model. That is, at optimality  $k \leq K$ .

521

522

### 523 3.2.2. Step 2: Component sharing through X-means clustering

524

525 If a feasible design solution exists after Step 1, X-means data mining clustering technique is once  
526 again employed, this time to determine the most similar product architecture design solutions  
527 within the product portfolio. While the first X-means clustering technique helped identify the  
528 similar groups of attributes in the raw data, the X-means clustering employed in Step 2 will  
529 help identify the similar groups of design variable values among the feasible product architecture  
530 design solutions (Step 2 in Figure 1).

531 For a given vector of design variables  $(\mathbf{x}_k)$  of an optimal product architecture solution, (where  
532 the objective function  $F(\mathbf{x}_k)$  of product architecture ( $k$ ) has been minimized given the external  
533 targets  $\mathbf{T}^{C_j}$  and the local objective(s)  $f_k(\mathbf{x}_k)$ ), the goal is to determine the similarity among product  
534 architecture variable solutions. The notion is that the closer the optimal design solutions are, for  
535 example  $[(\mathbf{x}_k)$  and  $(\mathbf{x}_{k+1})]$ , the more likely these product architectures may be able to share certain  
536 design components.

537

538

### 539 3.2.3. Step 3: Product family optimization with shared design components

540

541 The third and final step in the proposed product family design methodology aims to reduce  
542 the product portfolio cost by sharing certain components among product architectures, thereby  
543 creating a family of products (Step 3 in Figure 1). Since the component sharing decision is  
544 inherently a combinatorial problem, Step 2 of the design methodology eliminates the need to search  
545 all possible component sharing combinations by guiding the component sharing decisions based  
546 on the optimal solution of each resulting product architecture. Once similar product architectures  
547 have been identified by the X-means technique in Step 2, the component variables are identified  
548 and modelled as linking variables  $(\mathbf{y}_{s,k})$  in the quasi-separable bi-level problem. The model in  
549 Step 1 is adapted into a bi-level hierarchical optimization model where *level 1* strictly handles  
550 the coordination of the linking variables and *level 2* still remains the product architecture design  
level, but this time including the linking variable targets as part of the objective function. The

551 bi-level design problem is modelled based on the quasi-separable problem (Kim 2001, Kim *et al.*  
 552 2002, Kokkolaras *et al.* 2002, Allison *et al.* 2006). A bi-level model is presented which comprises  
 553 the component sharing co-ordination model at the upper level and the individual product design  
 554 model at the lower level. At the component sharing level, updated linking variable values are  
 555 distributed among product variants in an iterative manner until a feasible solution is achieved that  
 556 is common among all product variants. If a feasible design solution does not exist for a given  
 557 sharing scenario (that is, linking variable value  $\mathbf{y}_s$  does not converge to a solution shared by all  
 558 products), the original product design solutions (without shared variables) from Step 1 are kept.

559

560

561

*Upper level: Component sharing co-ordination*

562

563

564

565

566

567

568

569

570

571

572

Minimize

573

574

$$\varepsilon_y \quad (8)$$

575

576

Subject to:

577

578

579

$$g1 : \sum_{k \in Q} \left\| \mathbf{y}_s - \mathbf{y}_{s,k}^{Eng} \right\|_2^2 - \varepsilon_y \leq 0 \quad (9)$$

580

Here,

581

582

583

584

585

$\mathbf{y}_s$  Linking variable at the upper level. In essence,  $\mathbf{y}_s$  is simply a coordination variable ensuring that at the optimal solution, all of the subsystems attain the same value. Equation (9) is always active in the above formulation so solving for  $\mathbf{y}_s$ , it can be observed that at each iteration  $\mathbf{y}_s$  assumes the average value of the linking variable(s) being shared across the products within the product family.

586

587

588

$\mathbf{y}_{s,k}^{Eng}$  Linking variable value at the lower level cascaded to the upper level. This is constant at each iteration in the above formulation that is subsequently updated at the engineering product architecture optimization level after each iteration.

589

590

591

592

593

594

595

596

$k$  The  $k$ th candidate product architecture that has been identified for component sharing.

$Q$  The total number of products that exist in a particular candidate product family. This is based on the X-means cluster solutions described in Step 2. The term candidate product family is used because until a feasible design solution can be achieved for the shared component case, these  $Q$  products will remain unique products within the product portfolio (note that  $Q \leq K$  which simply means that the number of candidate product families cannot exceed the total number of unique products that initially exist).

597

598

599

600

$\varepsilon_y$  Deviation tolerance between linking variables. For each shared variable, another constraint  $g(i)$  is added based on a similar formulation as equation (9) and add another tolerance variable in the objective function to represent this additional shared variable.

601 *Lower level: Product family optimization*

602

603 In the  $k$ th sub-problem,

604 Minimize

$$605 \quad F(x)_{Architecture(k)} = f_k + \mathbf{w}' \left\| \mathbf{T}^{C_j} - \mathbf{R}_k^{Eng} \right\|_2^2 + \left\| \mathbf{y}_s^U - \mathbf{y}_{s,k} \right\|_2^2 \quad (10)$$

606

607 Subject to:

608

609

610

611

612

613  $f_k$  Local product design objective function (s).

614  $\mathbf{T}^{C_j}$  Vector of product attributes represented by the cluster centroid in the data mining  
615 model. That is, for cluster centroid  $C_j = [A_1, A_2, \dots, A_p]$ , target  $\mathbf{T}^{C_j}$  is set as  
616  $[A_1, A_2, \dots, A_p]$  where  $A_1, A_2, \dots, A_p$  represent attribute values for a given  
617 centroid  $C_j$ .

618  $\mathbf{w}'$  The vector of newly transformed RELIEFF weights for target vector  $\mathbf{T}^{C_j}$ .

619  $\mathbf{R}_k^{Eng}$  Vector of engineering responses based on local design variables.  $\mathbf{R}_k^{Eng}$  is a function of  
620 local design variables  $\mathbf{x}_k$ , and is represented by  $\mathbf{R}_k^{Eng}(\mathbf{x}_k, \mathbf{y}_{s,k})$ .

621  $\mathbf{g}_k$  Inequality design constraints.

622  $\mathbf{h}_k$  Equality design constraints.

623  $\mathbf{y}_s^U$  Linking variable target value cascaded down to the lower level from the upper level; a  
624 constant value at each iteration that is subsequently updated with each successful  
625 iteration.

626  $\mathbf{y}_{s,k}$  Linking variable at the lower level. This is local to the  $k$ th model and attempts to match  
627 the value of  $\mathbf{y}_s^U$  at each iteration.

628

629

The overall flow of the proposed product family optimization is succinctly described below:

630

### ***Bi-level product family optimization***

631

#### **Step 1:**

632 Given  $\mathbf{w}'$  vector of weights and  $\mathbf{T}^{C_j}$  targets, where  $\text{length}(\mathbf{w}') = \text{length}(\mathbf{T}^{C_j})$  and  $K$  cluster  
633 centroids:

634

635 1. Solve  $K$  engineering design problems (**with no** linking variables  $\mathbf{y}_{s,k}$ ) weighting each  
636  $\left\| \mathbf{T}^{C_j} - \mathbf{R}_k^{Eng} \right\|_2^2$  based on RELIEFF;

637 2. If solution exists for the Individual Product Design Optimization Model (i.e., *optimal*  
638  $\left\| \mathbf{T}^{C_j} - \mathbf{R}_k^{Eng} \right\|_2^2$  solution while satisfying local objectives and constraints);

639 3. Optimal solution found for weights  $\mathbf{w}'$  and targets  $\mathbf{T}^{C_j}$  without sharing components;

640

641

#### **Step 2:**

642 4. Employ **X-means** clustering to identify candidate product families based on solution  
643 similarities from **Step 1**;

644

645

#### **Step 3:**

646 5. Solve bi-level quasi-separable problem (component sharing among products) using the  
647 *Upper Level-Lower Level* formulation **with** linking variables  $\mathbf{y}_{s,k}$ ;

648 6. If feasible solution exists (i.e., *optimal*  $\left\| \mathbf{T}^{C_j} - \mathbf{R}_k^{Eng} \right\|_2^2$  and  $\left\| \mathbf{y}_s - \mathbf{y}_{s,k}^L \right\|_2^2$  at the *Lower Level*  
649 and also *optimal*  $\varepsilon_y$  at the *Upper Level*, ( $\varepsilon_y$  should be close to 0 at the *Upper Level*, indicating  
650 a feasible shared component among product variants within a product family));

- 651 7. Optimal solution found for weights  $w'$  and targets  $T^{Cj}$  and linking variables  $y_{s,k}$  for each
- 652 product variant;
- 653 8. Else, solution does not exist for linking variable scenario; that is, sharing  $y_{s,k}$  is not feasible
- 654 for product variants, therefore keep initial solutions found from **Step 1**;
- 655 9. end;
- 656 10. end;

#### 660 4. Application: Aerodynamic particle separator case study

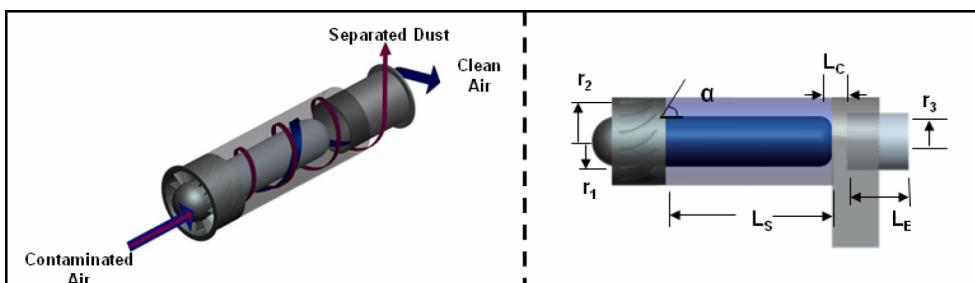
661  
662 Indoor air quality (IAQ) is becoming an increasing concern for human health. Particulate matter is  
663 a leading cause of human respiratory illness in addition to degrading the performance of heating  
664 ventilation and air conditioning (HVAC) systems (WHO). As a result, abatement technologies  
665 for these aerosols are in high demand. Aerodynamic particle separators are filter-less air clean-  
666 ing devices that can be capable of removing micron size particles with low energy consumption  
667 and minimal maintenance (Zhang 2005). Determining the optimal aerodynamic particle separa-  
668 tor design for a specific application is challenging when taking into account its unique system  
669 requirements and environmental conditions.

#### 671 4.1. The engineering design problem

##### 673 4.1.1. Aerodynamic particle separator design

674  
675 To demonstrate the effectiveness of the proposed design methodology, the design of a uniflow type  
676 particle separator is investigated (illustrated in Figure 3). The basic design of this device can be  
677 partitioned into three sections: (1) vane section, (2) straight region and (3) converging region/dust  
678 bunker. These sections are defined by eight design variables as shown in Figure 3 and Table 1.

679 The performance of an aerodynamic particle separator design is strongly dependent on system  
680 requirements and environmental conditions. System requirements such as the air cleaning effi-  
681 ciency, pressure drop (thus power consumption), air flow rate and overall device size contribute  
682 to the design objectives and directly define the constraints for a given application. Environmen-  
683 tal conditions, including the air properties and contaminant particle size distribution can have a  
684 significant impact on the performance of a particular design and must also be incorporated into  
685 the system model (Barker 2008). Together, these two groups can be used to characterize a given  
686 application or operating state. These factors can vary by an order of magnitude between different  
687 applications, thereby complicating the design process. The product architecture design objective  
688 will, therefore, be to minimize cost while satisfying external product preference targets and local  
689



690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700 Figure 3. Uniflow type aerodynamic particle separator flow pattern and design variables.

Table 1. Design variable notation for aerodynamic particle separator.

Variable	Units	Description
$r_1$	Meters (m)	Inner tube radius
$r_2$	Meters (m)	Outer tube radius
$L_S$	Meters (m)	Maximum pressure drop
$\alpha$	Radians (rad)	Vane discharge angle
$L_C$	Meters (m)	Length of converging gap
$r_3$	Meters (m)	Radius of exit tube
$L_E$	Meters (m)	Length of exit tube
$N$	#	Number of units in parallel

design constraints. The broader enterprise portfolio objective will be to minimize overall product design and development costs by capturing the component sharing opportunities that exists within the product portfolio.

## 4.2. Data mining product preferences

### 4.2.1. Raw data set of product operating states

A data set of 1000 operating states was generated to simulate the large variation in physical requirements and environmental conditions characterizing the broad range of applications in which aerodynamic particle separators are frequently employed (Barker 2008). Table 2 represents a snapshot of the 1000 operating states with distinct product attributes and environmental conditions represented by each column. Section 5 of the article presents the results from both the RELIEFF attribute ranking algorithm and the X-means data mining clustering approach and demonstrates how the data mining process influences product family design efforts.

### 4.2.2. RELIEFF attribute weighting

The results from the RELIEFF attribute ranking approach in Table 3 reveal that the two engineering design targets—efficiency ( $\zeta$ ) and flow area ( $AF_{\max}$ )—have normalized weights of 0.1687 and 0.0785 respectively. The other attributes in Table 3 are used as design parameters in the engineering design model and, therefore, also play a significant role in the overall optimal solution. The efficiency ( $\zeta$ ) and flow area ( $AF_{\max}$ ) are selected based on the type of engineering problem being solved (in other applications, one may choose to set all attributes in the data set as targets for the product architecture design model). This vital information is a data pre-processing step that will help generate product cluster centroids that take into account the weighted attribute preferences of the different operating states given by the raw data set.

*Note:* RELIEFF results were obtained using Weka version 3.5.8 (Witten and Frank 2005) and it took approximately 20 seconds running on a Intel Pentium Duo 2.5 GHz Processor. The normalized attribute weights are based on Equation (4) in Section 3.1.1.

Table 2. Snapshot of aerodynamic particle separator data set consisting of 1000 states.

	Q (m <sup>3</sup> /s)	$\Delta P_{\max}$ (Pa)	$L_{\max}$ (m)	$AF_{\max}$ (m <sup>2</sup> )	$N_{\max}$ # units	F( $d_p$ ) (%)	$\rho_p$ (kg/m <sup>3</sup> )	$T_{\text{air}}$ (°C)	$P_{\text{air}}$ (kPa)	Efficiency %	Price \$
State 1	1.58	250	1	0.5	50	A1	2650	20	101	82	1,200
State 2	1.48	200	0.3	0.1	3	A4	2650	0	99	85	450
State 3	1.27	1500	1.5	1.5	16	Limestone	2700	500	200	90	900

Table 3. Attribute ranking of raw dataset via RELIEFF algorithm.

Attribute Rank	Attribute Name	Attribute Weight	Normalized Attribute Weight
Highest	Nmax	0.0881	1.0000
	Efficiency	0.0110	0.1687
	Tair	0.0100	0.1575
	Pair	0.0099	0.1570
	Q	0.0073	0.1287
	Dpmax	0.0036	0.0888
	Afmax	0.0026	0.0785
	Lmax	0.0018	0.0695
	Rhop	-0.0006	0.0433
	Operating States	-0.0031	0.0169
Lowest	Fdp	-0.0046	0.0000

#### 4.2.3. Data mining X-means clustering results

The X-means clustering results reveal that a total of five clusters most accurately represents the similarities in the data set of 1000 operating states. The results from Table 3 represent the product design targets and parameters for the product portfolio of aerodynamic particle separators. Initially, each product centroid will be used to design an individual aerodynamic particle separator. Component sharing benefits will then be presented based on the vane section component. [Results attained using Weka version 3.5.8 (Witten and Frank 2005) and Data to Knowledge D2K (McEntire 2003)].

### 4.3. Engineering design optimization of product family

#### 4.3.1. Step 1: Individual product design optimization

The X-means clustering algorithm generates  $k = 1, \dots, 5$  clusters, each with unique centroids  $C_j$ . Based on the results from the X-means clustering, and the RELIEFF attribute weights accompanying each cluster centroid, engineers can now determine whether an optimal product design solution exists based on the aerodynamic particle separator response model.

The aerodynamic particle separator objective function attempts to match the particle separation efficiency target ( $\zeta_k^{C_j}$ ) and the flow area target ( $AF^{C_j}$ ) generated from the X-means clustering results while at the same time minimizing product design and manufacturing cost objective. The attributes within a cluster centroid ( $C_j$ ) will form the design/environmental parameters of the model. The efficiency model selected was initially developed by Zhang (2005) and later augmented by Barker (2008). In this model, the flow is assumed to be fully turbulent and the steady state particle motion results from a balance between the centrifugal force and aerodynamic drag in the Stokes regime (Zhang 2005). The vector  $\mathbf{x}$  contains the eight design variables as described by Table 1 and Figure 3. The cost function was based on the estimated mass of material required and injection moulding costs of the vane section. The material selected is an engineered polymer with a density of  $1200 \text{ kg/m}^3$  at a cost of \$3.00 per kilogram. The injection moulding cost is estimated at a fixed cost of \$10,000 per design for the required capital equipment and labour. The efficiency model as a function of variables in  $\mathbf{x}$  and particle size  $d_{pi}$  is shown in Equation (11). The total efficiency for a given particle size distribution is then calculated by Equation (12).

$$\xi(\mathbf{x}, d_{pi}) = 1 - \exp\left(-\frac{\rho_p d_{pi}^2 C_c Q \tan(\alpha) L_S}{9\eta(r_2^2 - r_1^2)}\right) \cdot \exp\left(\frac{\rho_p d_{pi}^2 C_c (V_t^2 G_t(\mathbf{x}) + V_z^2 G_r(\mathbf{x}))}{\eta V_z}\right) \quad (11)$$

$$\xi_T = \sum_{i=1}^N \xi(\mathbf{x}, d_{p_i}) \cdot F(d_{p_i}) \quad (12)$$

804 Here,

- 805  
 806  $Cc$  Cunningham slip correction factor.  
 807  $d_{p_i}$  Diameter of particle ( $i$ ), m.  
 808  $F(d_p)$  Particle size distribution.  
 809  $G_t(x)$  Efficiency model geometric relationship between design variables, tangential  
 810 acceleration.  
 811  $G_r(x)$  Efficiency model geometric relationship between design variables, radial acceleration.  
 812  $\rho_p$  Particle density, kg/m<sup>3</sup>.  
 813  $\eta$  Air viscosity, Pa·s or kg· m/s.  
 814  $Q$  Air flow rate, m<sup>3</sup>/s.  
 815  $V_t$  Tangential velocity of particle mixture.  
 816  $V_z$  Axial velocity of particle mixture.  
 817  $r_1$  Inner tube radius.  
 818  $r_2$  Inner tube radius.  
 819  $\alpha$  Vane discharge angle.  
 820  $L_S$  Maximum pressure drop.

821 The engineering design model for the aerodynamic particle separator can be mathematically  
 822 represented as:

823 *k*th aerodynamic particle separator

824 Minimize:

$$F(x)_{Architecture(k)} = w'_\zeta \left\| \zeta_k^{Cj} - \zeta_k^{Eng} \right\|_2^2 + w'_{AF} \left\| AF_k^{Cj} - AF_k^{Eng} \right\|_2^2 + Cost_k \quad (13)$$

828 Subject to:

829 Pressure drop constraint (g1):

$$P_T(\mathbf{x}) - P_{\max} \leq 0 \quad (14)$$

832 Face area constraint (g2):

$$4r_2^2 N - AF_{\max} \leq 0 \quad (15)$$

834 Product length constraint (g3):

$$L_V + L_S + L_C + L_E - L_{\max} \leq 0 \quad (16)$$

837 Here,

- 839  $AF_k$  Maximum allowable face area perpendicular to air flow direction.  
 840  $w'_\zeta$  Efficiency RELIEFF attribute weight.  
 841  $w'_{AF}$  Flow area (AF) RELIEFF attribute weight.  
 842  $L_{\max}$  Total allowable length of the system.  
 843  $L_V$  Length of vane section.  
 844  $L_S$  Length of straight region.  
 845  $L_C$  Length of converging region.  
 846  $L_E$  Length of exit tube.  
 847  $P_T(\mathbf{x})$  Total pressure drop of the system as a function of design variables  $\mathbf{x}$ .  
 848  $N$  Number of aerodynamic particle separator units in one module.  
 849  $P_{\max}$  Maximum allowable pressure drop (air flow restriction).  
 850  $Cost_k$  Total product cost represented as the summation of individual component costs.

851 *Note:* The design model is also bounded by a set of linear inequality constraints  $Ax \leq b$  and con-  
 852 straints Equations (14)–(16) that can be further expanded. A more detailed design model can be found  
 853 in Barker (2008).

#### 854 855 4.3.2. Step 2: Component sharing through X-means clustering

856  
857 If an optimal solution exists for the aerodynamic product portfolio based on the X-means clustering  
 858 targets, the next step is to determine whether additional costs savings can be realized by shar-  
 859 ing the most design intensive components among different product architectures. The X-means  
 860 clustering technique is employed to determine the similarities among the unique aerodynamic  
 861 particle separator designs based on the solution results after Step 1. A successful sharing solution  
 862 among products represents a unique product family. The results from the unique aerodynamic  
 863 particle separator solutions can be seen in Table 5 which is further explained in Section 5.1.

#### 864 865 4.3.3. Step 3: Product family optimization with shared design components

866  
867 For the aerodynamic particle separator case study, the vane section is the most design intensive  
 868 and costly component. The complex curved vanes must be injection moulded, which requires a  
 869 unique mould to be machined for each vane section design. By employing the X-means clustering  
 870 technique, product engineers will be able to (1) determine which product architecture designs are  
 871 similar based on the solutions attained during Step 1 and (2) determine the number of candidate  
 872 product families to include in the enterprise product portfolio based on the number of X-means  
 873 cluster centroids generated. The L2 norm distance measure used by X-means will favour those  
 874 design solutions that are numerically close to one another. This will help guide the sharing decision  
 875 of the vane section as products with close numerical values for the variables that define the vane  
 876 section (vane angle  $\alpha$ , the inner and outer tube radii  $r_1$  and  $r_2$ ) will be favoured within a given  
 877 cluster centroid.

#### 878 879 *Upper level: Component sharing co-ordination*

880  
881 The upper level (component sharing co-ordination) of the aerodynamic particle separator model  
 882 will handle the co-ordination of the shared vane section among product families. The component  
 883  
884  
885

886 Table 4. Product cluster centroids based on X-means clustering algorithm.

States	System requirements					System requirements					
	Q (m <sup>3</sup> /s)	$\Delta P_{\max}$ (Pa)	$L_{\max}$ (m)	$A_{F\max}$ (m <sup>2</sup> )	$N_{\max}$ # units	F( $d_p$ ) (%)	$\rho_p$ (kg/m <sup>3</sup> )	$T_{\text{air}}$ (°C)	$P_{\text{air}}$ (kPa)	Efficiency %	Price \$
<i>Product Centroid 1</i>											
285	1.20	463	0.71	0.60	36	Limestone	2226	25	100	78	962
<i>Product Centroid 2</i>											
765	3.39	1507	0.79	0.68	39	A4	2211	71	111	86	1,168
<i>Product Centroid 3</i>											
750	3.36	1623	0.67	0.55	14	Limestone	2389	72	107	85	411
<i>Product Centroid 4</i>											
456	1.94	911	0.72	0.62	24	A1	2278	45	102	80	632
<i>Product Centroid 5</i>											
260	1.29	458	0.69	0.75	11	Limestone	2225	24	100	78	290

901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950

Table 5. Optimal solutions for individual aerodynamic particle separator designs.

Product	Design Variables										Product Efficiency $\zeta$ (%)	Product Unit Cost \$	# Units/Cluster (# units)	Injection Mold Cost \$	Total Product Cost \$
	$r_1$ (meters)	$r_2$ (meters)	$L_s$ (meters)	$\hat{a}$ (rad)	$L_c$ (meters)	$r_3$ (meters)	AF (meters)	L (meters)	N (# unit)						
Particle Separator 1	0.1732	0.1936	0.1098	1.0400	0.0109	0.1163	0.1207	0.6000	4	77.9994	69.01	226	10,000.00	25,595.41	
Particle Separator 2	0.0707	0.0842	0.2709	1.0400	0.0088	0.0715	0.27971	0.6799	24	85.9915	171.03	207	10,000.00	45,403.45	
Particle Separator 3	0.0248	0.0991	0.3539	1.0400	0.0991	0.0527	0.45301	0.5499	14	84.9963	865.53	173	10,000.00	24,970.02	
Particle Separator 4	0.0691	0.0804	0.3696	1.0400	0.0804	0.0683	0.44998	0.6199	24	79.9955	189.70	233	10,000.00	54,199.65	
Particle Separator 5	0.1758	0.1936	0.0932	1.0400	0.0095	0.1068	0.10267	0.7500	5	77.9993	85.36	161	10,000.00	23,742.41	
Total Product Portfolio cost														173,910.93	

951 sharing objective function will minimize the tolerance deviation variable of each shared variable.  
 952 There are three variables that define the vane section, including the vane angle  $\alpha$ , the inner and  
 953 outer tube radii  $r_1$  and  $r_2$ .

954 **Minimize**

$$955 \quad \varepsilon_\alpha + \varepsilon_{r_1} + \varepsilon_{r_2} \quad (17)$$

957 **Subject to:**

$$958 \quad g1 : \left\| \alpha_s - \alpha_{s,k}^{Eng} \right\|_2^2 - \varepsilon_\alpha \leq 0 \quad (18)$$

$$959 \quad g2 : \left\| r_{1,s} - r_{1,s,k}^{Eng} \right\|_2^2 - \varepsilon_{r_1} \leq 0 \quad (19)$$

$$960 \quad g3 : \left\| r_{2,s} - r_{2,s,k}^{Eng} \right\|_2^2 - \varepsilon_{r_2} \leq 0 \quad (20)$$

966 Here,

- 967  $\alpha_s$  Vane angle linking variable at the component sharing level.  
 968  $\alpha_{s,k}^{Eng}$  Value of vane angle linking variable response of engineering design model for product  
 969  $k$ .  
 970  $r_{1,s,k}$  Inner tube radius ( $r_1$ ) linking variable at the component sharing level.  
 971  $r_{1,s,k}^{Eng}$  Value of inner tube radius ( $r_1$ ) linking variable response of engineering design model  
 972 for product  $k$ .  
 973  $r_{2,s,k}$  Outer tube radius ( $r_2$ ) linking variable at the component sharing level.  
 974  $r_{2,s,k}^{Eng}$  Value of outer tube radius ( $r_2$ ) linking variable response of engineering design model  
 975 for product  $k$ .  
 976  $\varepsilon_\alpha$  Deviation tolerance variable between vane angle linking variable that is minimized in  
 977 the objective function.  
 978  $\varepsilon_{r_1}$  Deviation tolerance variable between inner radius linking variable that is minimized in  
 979 the objective function.  
 980  $\varepsilon_{r_2}$  Deviation tolerance variable between outer radius linking variable that is minimized in  
 981 the objective function.  
 982  
 983  
 984

985 To minimize overall product portfolio costs, the number of unique vane section designs will  
 986 be minimized by sharing this component with products that can attain a feasible design solution  
 987 given this added objective. Equation (13) is, therefore, reformulated to reflect the candidate product  
 988 families and also the shared vane components among each of these products within a given product  
 989 family (represented as linking variables).

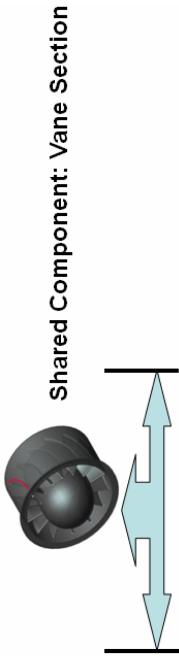
990 *Lower level: Product family optimization*

991 **Minimize:**

$$992 \quad F(x)_{Architecture(k)} = w'_\zeta \left\| \zeta_k^{C_j} - \zeta_k^{Eng} \right\|_2^2 + w'_{AF} \left\| AF_k^{C_j} - AF_k^{Eng} \right\|_2^2 + Cost_k + \left\| \alpha_{s,i} - \alpha_{s,k}^{Link} \right\|_2^2 \\
 993 \quad + \left\| r_{1s} - r_{1s,k}^{Link} \right\|_2^2 + \left\| r_{2s} - r_{2s,k}^{Link} \right\|_2^2 \quad (21)$$

994 **Subject to:** Constraints as defined in Equations (14), (15) and (16).  
 1000

Table 6. Optimal solutions for aerodynamic particle separator product families sharing the vane component.\*



Product	Product Design Variables										Product Efficiency	Product Unit Cost	# Units/Cluster	Injection Mold Cost	Total Product Cost
	$r_1$ (meters)	$r_2$ (meters)	$L_s$ (meters)	$\acute{\alpha}$ (rad)	$L_c$ (meters)	$r_3$ (meters)	AF (meters)	L (meters)	N (# unit)	$\zeta$ (%)					
<b>Particle Separator Product Family 1</b>															
Variant 1	0.1643	0.1935	0.2771	0.9532	0.0106	0.0884	0.5993	0.2877	4	77.9985	81.54	10,000.00	28,428.21		
Variant 5	0.1653	0.1935	0.2991	0.9533	0.0097	0.0812	0.7487	0.3088	5	77.9979	105.10	Shared	16,921.90		
<b>Particle Separator Product Family 2</b>															
Variant 2	0.0908	0.1064	0.5559	0.9446	0.0071	0.0905	0.6797	0.5630	15	95.84444	196.45	10,000.00	50,664.52		
Variant 3	0.0903	0.1070	0.0509	0.9446	0.0072	0.0686	0.5495	0.0581	12	84.9991	74.77	Shared	12,935.74		
<b>Particle Separator Product Family 3</b>															
Variant 4	0.069059	0.080357	0.369626	1.04	0.080357	0.068303	0.449983	0.619893	24	79.995506	189.70	10,000.00	54,199.65		
<b>Total Product Portfolio cost</b>												163,150.02			

\*Optimal results attained using Matlab® and Tomlab® to solve the mixed integer nonlinear programming problem (Griffiths 2005, Holmstrom *et al.* 2006).

## 5. Results and discussion

### 5.1. Aerodynamic particle separator optimization results

Given the product design targets from the data mining X-means clustering step, the aerodynamic particle separator model first attempts to identify feasible design solutions for the efficiency ( $\zeta^{C_j}$ ), flow area ( $AF^{C_j}$ ) targets and given physical and environmental ( $T_{air}$ ,  $P_{air}$ , etc.) parameters for each unique cluster centroid ( $C_j$ ). The aerodynamic particle separator solutions in Table 5 reveal that a total of five unique products can be designed for the initial five cluster centroids targets generated by the X-means clustering with a total product portfolio cost of \$173,910.

If Step 1 of the product family design methodology is successful, engineers can further investigate the potential costs savings (Steps 2 and 3) that may be realized due to component sharing. The X-means clustering technique performed during Step 2 reveals that out of the five unique aerodynamic particle separator solutions, products 1 and 5 form a feasible unique product family cluster, products 2 and 3 another and finally product 4 cannot be shared with any other product and, therefore, reverts back to the original solution from Step 1. The cost of the injection mould manufacturing process presents an opportunity for the initial product portfolio of five unique products to be redesigned. The vane section of the product which is made through the injection moulding process is shared among similar products existing in the original portfolio. In this case study, the decision to share the vane angle is known a priori due to the high cost of designing each individual injection mould for the vane. Step 3 of the product family design methodology employs the X-means clustering algorithm to identify products that have similar vane design solutions. The decision to share the vane angle is an attempt to minimize the overall costs of the enterprise product portfolio by minimizing the number of unique vane sections needed for the five aerodynamic particle separators. Products successfully sharing a vane section will be considered a unique product family and each product existing in this product family, is defined as a variant. However, it must be noted that the cost savings benefits of component sharing using the product family approach to design may be offset by the decrease in the performance capabilities attainable by the newly designed product variants. This trade-off scenario will, therefore, be based on how much cost savings can be realized through component sharing and how much performance deviation can be accommodated by the customer.

The results in Table 6 reveal that sharing products 1 and 5, 2 and 3 (with product 4 being a separate unique design), reduces the total product portfolio cost to \$163,150; a total savings of approximately \$10,760 for this product portfolio design scenario. However, it can be observed that the efficiency of product 2 decreases from 85.99% with the individual optimization model solution (Table 5) to 85.84% with the component sharing product family model solution (Table 6). The level of allowable performance deviation will be dependent on customer expectations and the level of competition within the market space. Although a feasible design may not always exist for every sharing scenario (for example sharing a single vane component for each of the five products returned an infeasible solution), the benefits of investigating sharing strategies through the X-means clustering recommendations may prove beneficial as can be seen from the results in Table 6.

## 6. Conclusion

In this work, a comprehensive product family design methodology is presented that integrates realistic product operation data with the engineering design of complex products such as the aerodynamic particle separator. The data mining RELIEFF algorithm is employed to determine the weights of each attribute. This information is then incorporated into the data mining X-means clustering algorithm in order to generate the number of clusters along with the cluster centroids

1101 that are inherent to the data itself. The results of the data mining clustering technique aid in  
 1102 determining the number of unique products to design for a group of highly diverse customers.  
 1103 With this clustering information, a product architecture can be designed that takes into account  
 1104 specific customer product functionality needs that are represented in a large data set. Further cost  
 1105 savings can be realized through a component sharing strategy that is achieved in this work by  
 1106 once again employing the X-means clustering technique to identify similar design solutions. The  
 1107 hope is to expand on the concepts presented in this work by enabling the feasibility of the product  
 1108 architecture optimization step to influence the generation of X-means cluster centroids. That is,  
 1109 local objective functions may be highly sensitive to certain local design variables which can be  
 1110 taken into account during the X-means clustering step.

1111

1112

### 1113 Acknowledgements

1114 This material is based on work supported by the National Science Foundation under Award No. 0726934. Any opinions,  
 1115 findings and conclusions or recommendations expressed in this publication are those of the authors and do not necessarily  
 1116 reflect the views of the National Science Foundation.

1117

1118

### 1119 Note

1120

- 1121 1. The initial version presented at the 12th AIAA/ISSMO Multidisciplinary Analysis and Optimization Conference,  
 1122 Victoria, British Columbia, Canada.

1123

### 1124 References

1125

1126 Agard, B. and Kusiak, A., 2004. Data-mining-based methodology for the design of product families. *International Journal*  
 1127 *of Production Research*, 42 (15), 2955–2969.

1128 Alexandrov, N.M. and Lewis, R.M., 2002. Analytical and computational aspects of collaborative optimization for  
 1129 multidisciplinary design. *AIAA Journal*, 40 (2), 301–309.

1129 Alizon, F., Khadke, K., Thevenot, H.J., Gershenson, J.K., Marion, T.J., Shooter, S.B. and Simpson, T.W., 2007. Frameworks  
 1130 for product family design and development. *Concurrent Engineering*, 15 (2), 187–199.

1131 Allison, J., Walsh, D., Kokkolaras, M., Papalambros, P.Y. and Cartmell, M., 2006. Analytical target cascading in aircraft **Q2**  
 1132 design. *44th AIAA aerospace sciences meeting and exhibit*, date, Reno, Nevada. Place: Publisher, 00–00.

1132 Arora, S., Raghavan, P. and Rao, S., 1998. Approximation schemes for euclidean k-median and related problems. **Q3**  
 1133 *Proceedings of the 30th annual ACM symposium theory of computing*, date, place. Place: Publisher, 106–113.

1134 Barker, D., 2008. *Development of an optimization design platform for aerodynamic particle separators*. Thesis (Masters).  
 1135 University of Illinois at Urbana-Champaign.

1135 Collier, D., 1981. The measurement and operating benefits of component part commonality. *Decision Sciences*, 12 (1),  
 1136 85–96.

1136 Fayyad, U., Piatetsky-Shapiro, G. and Smyth, P., 1996. From data mining to knowledge discovery in databases. **Q4**  
 1137 *Communications of the ACS*, 39, n11, 24(3).

1138 Griffiths, D.F., 2005. *An introduction to matlab*. Nethergate, Dundee, DD1 4HN, Scotland, UK: University of Dundee.

1139 Han, J. and Kamber, M., 2006. *Data mining concepts and techniques*. 2nd edn. Place: Morgan Kaufmann. **Q5**

1139 Hartigan, J.A. and Wong, M.A., 1979. A k-means clustering algorithm. *Applied Statistics*, 28 (1), 100–108.

1140 Holmstrom, K., Goran, A.O. and Edvall, M.M., 2006. *User's guide for tomlab/minlp*. 855 Beech St 121, San Diego, CA,  
 1141 USA: Tomlab Optimization Inc.

1142 Jain, A.K. and Dubes, R.C., 1988. *Algorithms for clustering data*. Place: Prentice Hall College. **Q6**

1143 Jiao, J. and Tseng, M.M., 2000. Understanding product family for mass customization by developing commonality indices.  
 1144 *Journal of Engineering Design*, 11 (3), 225–243.

1144 Kanungo, T., Mount, D.M., Netanyahu, N.S., Piatko, C.D., Silverman, R.S. and Wu, A.W., 2002. An efficient k-means  
 1145 clustering algorithm: Analysis and implementation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,  
 1146 24 (7), 881–892.

1147 Kass, R.E. and Wasserman, L., 1995. A reference Bayesian test for nested hypotheses and its relationship to the Schwarz  
 1148 criterion. *Journal of the American Statistical Association*, 90 (431), 928–934.

1148 Khajavirad, A. and Michalek, J.J., 2007. An extension of the commonality index for product family optimization. **Q7**  
 1149 *Proceedings of the ASME 2007 IDET and CIE conference*, date, Las Vegas, Nevada, USA. Place: Publisher,  
 00–00.

1150 Kim, H.M., 2001. *Target cascading in optimal system design*. Dissertation (PhD). The University of Michigan.

- 1151 Kim, H.M., Kokkolaras, M., Louca, L.S., Delagrammatikas, G.J., Michelena, N.F., Filipi, Z.S., Papalambros, P.L., Stein,  
1152 J.L. and Assanis, D.N., 2002. Target cascading in vehicle redesign: A class vi truck study. *International Journal of  
1153 Vehicle Design*, 3 (3), 199–225.
- 1154 Kim, H.M., Michelena, N.F. and Papalambros, P.Y., 2003. Target cascading in optimal system design. *Transactions of  
Q8 1155 ASME: Journal of Mechanical Design*, 125 (3), 474–480.
- 1156 Kira, K. and Rendell, L., 1992. A practical approach to feature selection. *Proceedings of international conference on  
1157 machine learning*, date, place. Aberdeen: Morgan Kaufmann, 249–256.
- 1158 Kokkolaras, M., Fellini, R., Kim, H.M., Michelena, N.F. and Papalambros, P.Y., 2002. Extension of the target cascading  
Q9 1159 formulation to the design of product families. *Structural and Multidisciplinary Optimization*, 24 (4), 293–301.
- 1160 Kononenko, I., 1994. Estimating attributes: Analysis and extensions of relief. *Proceedings of European conference in  
1161 machine learning*, date, place. Place: Publisher, 171–182.
- 1162 Kota, S., Sethuraman, K. and Miller, R., 2000. A metric for evaluating design commonality in product families. *Transactions  
1163 of ASME: Journal of Mechanical Design*, 122 (4), 403–410.
- 1164 Kusiak, A., 2006. Data mining: Manufacturing and service applications. *International Journal of Production Research*,  
1165 44 (19), 4175–4191.
- 1166 Martin, M. and Ishii, K., 1996. Design for variety: A methodology for understanding the costs of product proliferation.  
Q10 1167 *Proceedings of DETC 1997*, date, place. Place: Publisher, 00–00.
- 1168 Martin, M.V. and Ishii, K., 2002. Design for variety: Developing standardized and modularized product platform  
1169 architectures. *Research in Engineering Design*, 13 (4), 213–235.
- 1170 Martin, M.V. and Ishii, K., 1997. Design for variety: Development of complexity indices and design charts. *Proceedings  
1171 of DETC 1997*, date, place. Place: Publisher, 00–00.
- 1172 McAdams, D.A., Stone, R.B. and Wood, K.L., 1999. Functional interdependence and product similarity based on customer  
1173 needs. *Research in Engineering Design*, 11 (1), 1–19.
- 1174 McAdams, D.A. and Wood, K.L., 2002. A quantitative similarity metric for design-by-analogy. *Transactions of ASME:  
1175 Journal of Mechanical Design*, 124 (2), 173–182.
- 1176 Mcentire, J., 2003. *D2k toolkit user manual*. Place: Publisher.
- 1177 Messac, A., Martinez, M.P. and Simpson, T.W., 2002. Effective product family design using physical programming.  
Q13 1178 *Engineering Optimization*, 34 (3), 245–261.
- 1179 Moon, S.K., Kumara, S.R.T. and Simpson, T.W., 2006. Data mining and fuzzy clustering to support product family design.  
1180 *Proceedings of the ASME design automation conference*, date, place. Place: Publisher, 00–00.
- 1181 Pelleg, D. and Moore, A., 2000. X-means: Extending k-means with efficient estimation of the number of clusters.  
Q14 1182 *Proceedings of the 17th international conference on machine learning*, date, Stanford University, California. Place:  
1183 Publisher, 727–734.
- 1184 Siddique, Z., Rosen, D.W. and Wang, N., 1998. On the applicability of product variety design concepts to automotive  
1185 platform commonality. *Proceedings of DETC 1998*, date, Atlanta, GA. Place: Publisher, 00–00.
- 1186 Tarpey, T., 2007. A parametric k-means algorithm. *Computational Statistics*, 22, 71–89.
- 1187 Thevenot, H.J. and Simpson, T.W., 2007. A comprehensive metric for evaluating component commonality in a product  
1188 family. *Journal of Engineering Design*, 18 (6), 577–598.
- 1189 Tosserams, S., Etman, L.F.P. and Rooda, J.E., 2007. An augmented lagrangian decomposition method for quasi-separable  
1190 problems in mdo. *Structural and Multidisciplinary Optimization*, 34 (3), 211–227.
- 1191 Tucker, C.S. and Kim, H.M., 2007. Product family decision tree concept generation and validation through data mining  
1192 and multi-level optimization. *Proceedings of IDETC/CIE 2007*, date, Las Vegas, NV. Place: Publisher, 00–00.
- 1193 Tucker, C.S. and Kim, H.M., 2008. Optimal product portfolio formulation by merging predictive data mining with  
Q17 1184 multilevel optimization. *Transactions of ASME: ASME Journal of Mechanical Design*, 130 (4), 00–00.
- 1185 Wacker, J.G. and Trelevan, M., 1986. Component part standardization: An analysis of commonality sources and indices.  
1186 *Journal of Operations Management*, 6 (2), 219–244.
- 1187 WHO, [online]. Available from: [www.who.int/mediacentre/news/releases/2006/pr52/en/index.html](http://www.who.int/mediacentre/news/releases/2006/pr52/en/index.html) [Accessed September  
1188 2009].
- 1189 Witten, I.H. and Frank, E., 2005. *Data mining: Practical machine learning tools and techniques*. 2nd edn. Place: Morgan  
1190 Kaufmann.
- 1191 Zhang, Y., 2005. *Indoor air quality engineering*. Boca Raton, FL: CRC Press.