

Design as a Marked Point Process

John Quigley¹, Gokula Vasantha², Jonathan Corney³, David Purves¹, and Andrew Sherlock⁴

¹Department of Management Science, University of Strathclyde, G1 1XQ, Scotland

²School of Engineering and the Built Environment, Edinburgh Napier University, EH10 5DT, Scotland

³School of Engineering, University of Edinburgh, EH8 9YL, Scotland

⁴National Manufacturing Institute Scotland, University of Strathclyde, G1 1XJ, Scotland

Abstract

Although AI systems which support composition using predictive text are well established there are no analogous technologies for mechanical design. Motivated by the vision of a predictive system, that learns from previous designs and can interactively provide a list of established feature alternatives to the designer as a design progresses, this paper describes the theory, implementation and assessment of an intelligent system that learns from a family of previous designs and generates inferences using a form of spatial statistics.

The formalism presented, models 3D design activity as a ‘Marked Point Process’ that enables the probability of specific features being added at particular locations to be calculated. Because the resulting probabilities are updated every time a new feature is added the predictions will become more accurate as a design develops. This approach allows the cursor position on a CAD model to implicitly define a spatial focus for every query made to the statistical model. The authors describe the mathematics underlying a statistical model that amalgamates the frequency of occurrence of the features in the existing designs of a product family.

Having established the theoretical foundations of the work, a generic six step implementation process is described. This process is then illustrated for circular hole features using a statistical model generated from a dataset of hydraulic valves. The paper describes how the positions of each design’s extracted hole features can be homogenized through rotation and scaling. Results suggest that within generic part families (i.e. designs with common structure) a marked point process can be effective at predicting incremental steps in the development of new designs.

Keywords: Feature based Design, Predictive Design, Marked Point Process

1 Introduction

It has been argued that only 20% of design information is reused despite 90% of all design activities being based on the variants of existing designs [1], and that on average only 28% of design information is reused within manufacturing applications [2]. Design can be considered as a sequential decision-making process, where the current state of a design evolves through a series of design choices. A system is required where design features may be suggested to the designer for effective reuse, and these design reuse procedures can be learned from historical data [3, 4].

This paper introduces the underpinning mathematics required for implementation of a new generation of user interfaces that automatically identifies appropriate characteristics of previous designs for reuse based on a designer’s real time activity. As a design evolves the system generates predictions of the features which might be incorporated, and are informed by both previous work and the new, ongoing design. In order to identify the most relevant features, and avoid presenting the user with an overwhelming number of suggestions, the work reported exploits the location of information (i.e. features and mouse pointer) on a 3D Computer Aided Design (CAD) model so that predictions can be appropriate to specific positions on an engineering component. The system described assumes a single engineer developing a design by carrying out a series of operations on a CAD system. The system does not dictate any order of operations and allows the engineer’s focus to move around the component.

Designs seldom start with a blank sheet of paper, but are informed by past experiences with reports of as much as 75% of design activity comprising the re-use of existing knowledge [5]. In the context of designing industrial parts such activities comprise re-using, configuring, and assembling existing components. A key contributing factor to companies not performing projects on time and budget is the lack of knowledge re-use, which leads to frequent ‘reinventing the wheel’ rather than finding and using already known solutions [6].

Motivated by these observations this paper proposes a different form of design representation that can combine many design variations into a single probabilistic model that facilitates the reuse of previously used features during an interactive process that leads to the instantiation of a new design. By leveraging the available information, a probabilistic CAD system would prompt the engineers with fragments (i.e. features) of previously designed components to extend the current CAD design. Although reuse of common features in the design of many industrial products is desirable there could be cases where such a practice inhibits innovation. Aware of this the authors’ aim is not to automate but support with suggestions that the engineer is free to ignore. For this we propose modeling the design process as a Marked Point Process (MPP) to create a formal framework that can assess the association between designs. Similar approaches have been used successfully in neuroanatomy to analyze brain scan images through voxel based morphometry [7], as well a feature recognition in image analysis.

For our application, points are the coordinate location corresponding to where a design feature has

57 been placed and marks refer to the feature chosen. MPP is a form of spatial statistics, metrics based
58 on statistical tools that are used to characterize the distribution of events across space [8], and are
59 widely used across a number of application areas for example the distribution of trees in forests to
60 stars in the sky. Through this lens we view the behavior of engineers as a stochastic process, updating
61 throughout the design process on decisions made to place features in specific locations and thereby
62 supporting probabilistic measures for subsequent choices. The statistical inference can be supported
63 through historical data, viewing past designs as realizations from such a stochastic process. Specifically,
64 we develop a decision support system through a Bayesian methodology, where we start with a prior
65 distribution to assign a probabilistic measure on the features and location to be chosen by the engineer.
66 Following each choice, the prior distribution is updated to a posterior distribution based on this new data
67 and thereby making full use of all the information available. So as more design choices are made, the
68 model will be able to discriminate more effectively between historical designs based on similarity.

69 Given the above context and motivation, the authors defined the following goals for the work:

70 **Aim**

71 To define a computational framework that can support an interactive design process with suggestions of
72 features based on three inputs: a knowledge of existing designs; the state of an emerging design and a
73 location on the surface of the emerging design.

74 **Objectives**

- 75 1. Establish a method of homogenizing the orientation and dimensions of a collection of designs be-
76 longing to a product family.
- 77 2. Develop a statistical function that represents the probability of a particular feature occurring at a
78 particular location, on the surface of a design for a member of the product family.
- 79 3. Create a prototype implementation that can support an interactive design cycle which updates
80 the inferred probability of specific features occurring at given location as the design of a part is
81 modified.
- 82 4. Assess the accuracy of the feature predictions
- 83 5. Identify any inherent limitations or weakness in the approach

84 The rest of this paper is structured as follows: In Section 2 we provide a brief review of both predictive
85 design systems and relevant MPP literature to position the contribution of this work. In Section 3 we
86 present a generic overview of our process to support design development and in Section 4 we outline the
87 details of the mathematical model that underpins the process. We explore key characteristics of design

88 activities and data to inform our modeling choices, and we provide a generic modeling and decision
89 support framework. In Section 5, we evaluate the proposed modeling framework through a case study.
90 Finally, in Section 6 we reflect on the future direction of research in this area.

91 **2 Literature Review**

92 Systems for predictive design have to combine assessment of historical data with statistical methods so
93 that human users can easily choose, or ignore, suggestions that enhances the creative process. SMS text
94 messaging software, used by mobile phones, illustrate both the potential and challenges of engineering
95 useful predictive systems. However, while text prediction seek to identify patterns in a linear series of
96 symbols with a simple (i.e. keypad) interface, anticipating the intention of product designer requires
97 analysis of 3D information that has no canonical ordering (i.e. unlike a sentence of text, that reads from
98 left to right, a designer can essential edit shapes in any sequence). Despite these inherent difficulties
99 research into computational technologies that could enable predictive design systems has been reported
100 for more than a decade.

101 An early example is [9] who developed the “InspireMe” interface which allowed a user to ‘place’ and
102 ‘glue’ one of ten suggestions, proposed in response to a query shape, and then request new suggestions
103 for the resulting composite shape. The placed shape can be translated, rotated, and scaled to match the
104 query shape. The suggestions that are not useful can be removed and replaced with new suggestions.
105 [9] used a multi-dimensional histogram-based signature to encode shape’s global spatial structure and its
106 local detail to identify suggestions for a given shape query.

107 Later work recognized that there was potential to improve the accuracy of suggestions by combining
108 the frequency of occurrence with shape parameter values. For example, [10] demonstrated an interface for
109 an assembly-based modeling tool. The interface presents the user with semantical labeled tabs that can
110 be expanded hierarchically to show component sub-categories. The user can select a component and drag
111 it onto the current model. A probabilistic Bayesian network is then used to dynamically update both the
112 proposed component categories and the components based on their semantic and stylistic compatibility
113 with the current modeling state. The interface estimates whether the new component should have a
114 symmetric counterpart and computes the symmetry plane. Based on the modeling requirements, the
115 selected model can be moved to a position, rotated, scaled, duplicated, and glued.

116 While [9] focused on the design of assemblies of predefined component parts, [11] reported a predictive
117 system for component shape synthesis. Their approach provided an interactive platform for the user to
118 constrain shape synthesis based on high-level specifications (i.e. specific components, components from
119 particular categories, and components from learned latent styles) and an input shape database. Within
120 the proposed interactive shape synthesis interface a user can select constraints by selecting required: shape

121 styles, component categories and styles. The algorithm then proposes a list of synthesized objects based
122 on the given inputs. The discrete features help ensure that components selected for a synthesized shape
123 have compatible numbers of adjacent components of each type, and their edges have been identified and
124 stored with the category label of components (so they can be attached for placing where a component can
125 be attached to another component with symmetry relationships). Like [9], [11] also used a probabilistic
126 approach to identify and synthesize existing shapes from complex domains to generate new combinations
127 of components.

128 There is a tension in all reported work between accuracy and the number of predictions made. This
129 can be observed in [12] that describes a user interface that guides a designer’s selection with a list of
130 the 50 ‘best’ suggested components during an assembly based modeling process. The interface aims
131 to enable easy browsing and propose components that are most compatible with the current state of
132 the assembly design (represented as a 3D model). The interface allows users to manually drag, move,
133 scale, orient and combine selected components. The placed components can also be incorporated in the
134 design using Boolean operations (union, difference, and intersection) to obtain composite model. The
135 suggestion list automatically updates every time a component is added to the assembly. The suggestions
136 are ranked by size (larger components are given preference) at the start of the modeling process. The
137 marginal probability distribution computed from a factor graph by [12], which incorporates adjacency
138 and multiplicity factors of segmented components, to score and rank predicted components.

139 A different type of assembly design is considered by [13] who proposed an algorithm that takes a
140 partially completed 3D scene as input and propose relevant models in a user-specified region of interest
141 by leveraging text data. Suggestions are generated using three different approaches; Graph Kernel, N-
142 gram, and Merged. A query is generated by converting the given 3D scene into text that represents the
143 five closest models to a focal point nominated by the user. The algorithm uses co-occurrence, 5-gram
144 statistics from Google Web N-grams dataset and point-wise mutual information between the labels of
145 nearby models in the scene and the labels of models in the database to create suggestions.

146 For a very similar application [14] presented a method for generating novel arrangements of diverse 3D
147 objects synthesized from few given examples. The method creates a probabilistic model for scenes based
148 on Bayesian networks and Gaussian mixtures that can be trained by a small number of input examples
149 of relevant scenes retrieved from database. User were able to vary the degrees of similarity and diversity
150 in the generated scenes by controlling the weighting (through blending parameters) given to the influence
151 of the existing database of prior designs.

152 The “AttribIt” interface was developed in [15] which facilitates the targeted exploration of different
153 combinations of visual components using commands based on the relative semantic attribute. A user
154 initializes a design with a coherent combination of components from a database, then they select a subset
155 of these components and interactively increase, or decrease, the strength of an attribute using sliders.

156 In doing this they can observe changes to the whole design in real-time as new database components
 157 corresponding to the updated attribute strengths are swapped. The components are assembled automat-
 158 ically into a coherent design (provision for manual adjustments such as translation, rotation and scaling
 159 controls are available to refine the results). The interface shows regions of high geometric variation under
 160 the current attribute (highlighted in red color in Fig. 1).



Figure 1: User Interface for Assembly-based modeling using relative attributes [15]

161 When the overall form of a design (whether assembly or component) is constrained by function or the
 162 need to fit into a product family a template can be used to facilitate reuse. For example [16] developed
 163 templates that can be used in an interactive design system to create new 3D models in a design-by-
 164 example manner. The interface allows a user to choose template parts from the database, change their
 165 parameters, and combine them to create new models. The information in the template has been used
 166 to automatically position, align and connect parts by adjusting parameters, adding constraints, and
 167 assigning connectors. The assembly-based modeling system provides pick and drag substructures from
 168 different designs and add them to a working model. The elements on the selected node are represented
 169 in full color, while the others become semi-transparent during manipulation, and constrained degrees of
 170 freedom are hidden.

171 To support the generation of interior designs [17] developed a probabilistic hierarchical grammar
 172 to enable functional (rather than spatial) representation of an office environment. The aim was to
 173 support consistent segmentations, category labels and functional groupings of 3D scenes that characterizes
 174 geometric properties, cardinalities and spatial relationship in a hierarchical manner. A probabilistic
 175 grammar is used to automatically create consistent annotated scene graphs. Figure 2 illustrates an input
 176 scene mapped with labels and then converted into the hierarchical form using probabilistic grammar.
 177 A 7-dimensional descriptor (i.e. support and vertical relationships, horizontal separation, and overlap
 178 between objects) is used to describe the relationship between two objects. The dynamic programming
 179 for belief propagation was developed for scene parsing with optimal hierarchy. The technique creates
 180 candidate nodes based on spatial proximity, grammar binaries and finds the optimal binary hierarchy

181 which is converted to a logical hierarchy of the original grammar.

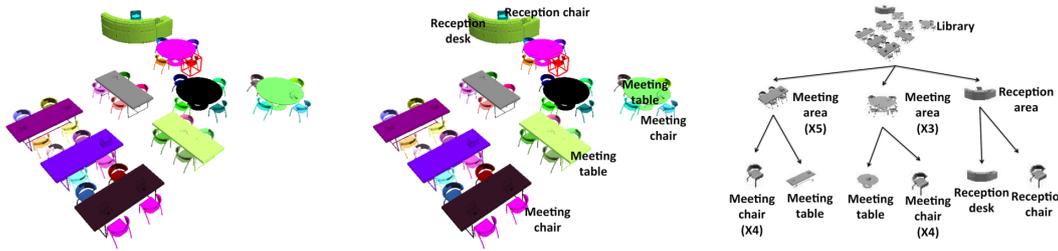


Figure 2: Input scene mapped with labels and converted into the hierarchical form using probabilistic grammar [17]

182 More recently [18] reported the “ComplementMe” user interface that aims to seamlessly integrate
183 suggested CAD models into the design process. A combination of embedding and retrieval neural network
184 architectures are proposed for suggesting complementary functional and stylistic components and their
185 placements within an incomplete 3D part assembly. The embedding network was used to map parts
186 to a low-dimensional feature space, and the retrieval network was used to retrieve partial assemblies to
187 appropriate components. The interface shows the possible candidates generated by sampling from the
188 conditional probability distribution predicted by the retrieval network. The user could select a desired
189 complementary component, and the algorithm predicts the location for it via the placement network. The
190 new shape will be synthesized for the user, and the next component is proposed based on the modified
191 assembly.

192 In conclusion the literature on predictive design systems is largely focused on the creation of assemblies
193 of 3D component models where frequently the positioning of suggested components is a manual task for
194 the user. In contrast the authors’ work is focused on the identification of shape features (i.e. fragments
195 of an entire model that are patterns of geometry such as holes) that are appropriate to a location defined
196 by the position of a user’s mouse pointer on the surface of a 3D object.

197 2.1 Marked Point Processes

198 Marked point processes (MPP) are widely applied within image analysis, where it was first introduced
199 by [19]. The methodology is used extensively and successfully for the extraction of multiple objects from
200 images. Applications include biological imagery on cells [20], disks in a plane [21], building outlines [22]
201 and person detection from camera images [23]. It is a flexible methodology that has been extended for
202 object extraction from images to arbitrarily shaped objects [24]. More recently, [25] have developed the
203 approach for microscope images, [26] have used MPP’s to automatically detect the locations of road
204 segments and [27] have used it for visual perceptions. A survey of marked point processes applied to
205 image analysis can be found in [28].

206 The literature to date has developed methods to extract images and characterize them in the form of
207 a MPP which are then stored in a database. Our focus complements this work, as we develop decision

208 support tools that also utilizes information about the location of extracted features in an MPP data
209 structure.

210 **3 Process for Constructing Marked Point Process Decision Sup-** 211 **port**

212 We propose a six step approach adapted from the CISSE process, see [29], for constructing empirical
213 prior distributions to support Bayesian analysis, which considers the following five steps; *Characterize*,
214 *Identify*, *Sentence*, *Select*, and *Estimate*. As described in the following for the third step we have placed
215 particular focus on homogenizing the data rather than sentencing the data and we have decomposed the
216 fifth step to consider prediction and updating.

217 Step 1: *Characterize the population of designs*. We begin by identifying those factors characterizing the
218 design. This is an important step because it defines the criteria by which data sets (i.e. historical
219 designs) are subsequently selected for inclusion in the comparator pool used to construct the
220 prior distribution. Examples of such characteristics may be with respect to types of layouts of
221 and/or features used within designs.

222 Step 2: *Identify candidate sample designs matching population*. The factors characterizing the population
223 of designs can be switched on/off for candidate designs effectively providing a means of making
224 a relative assessment of relevance against a set of criteria. We are simply trying to find the best
225 available data sets to make reasonable and timely inference. We are assuming that the current
226 design for which we are providing the decision support will be similar to one of these historical
227 designs. We can accommodate a unique apriori assessment on the likelihood of the current design
228 being realized to be like each possible candidate historical design, although our default may be a
229 uniform distribution prior.

230 Step 3: *Homogenizing the comparator data*. Generally, the higher the degree of homogeneity within the
231 comparator pool the more accurate the predictive inference [see 30]. This requires a measure for
232 similarity between designs, such as the KL divergence measure as proposed in [31] against which
233 the data can be transformed for homogeneity. Two key approaches to address this are scaling and
234 rotation. Firstly, all designs can be re-scaled into the unit cube. Secondly, the data describing
235 the locations of features can be rotated for alignment. This work should be performed prior to
236 the start of the design. This stage may be omitted if it is considered that information would be
237 lost in transforming the data, and the resulting prior would not be as effective at discriminating
238 between design types.

239 Step 4: *Select a probability model for the population of designs.* The nature of design patterns is such that
240 a parametric probability distribution is unlikely to exist that adequately represent the variability
241 of location and features within designs. As such, a non-parametric approach should be considered,
242 for which we recommend Kernel Density Estimation (KDE). Under such an approach, choices
243 will need to be made concerning the band-width parameter, which is essentially deciding on
244 allowable variation of location of features within similar designs. The resulting model is known
245 as the Feature Location Probability Function (FLPF), for which we would fit one to each historic
246 design to obtain a model for each design type.

247 Step 5: *Predictive Model.* The predictive distribution is simply a weighted average of the FLPF for each
248 design type in the comparator set, where the weights reflect the likelihood that the current design
249 will ultimately be realized as being similar to the candidate design in the set.

250 Step 6: *Update prior on design type and predictive distribution.* During the design process, Steps 5 and
251 6 are repeated in a cycle of feature addition and updating of the predictive distribution, which
252 we call the Predictive Feature Location Function (PFLF), to reflect how each change impacts on
253 the probable location of other features. This process, driven by the actions and selections of the
254 human designer, continues until the component part is complete (i.e. the design is finished).

255 Figure 3 provides a schematic for the predictive system. The data homogenization and FLPF can be
256 performed in advance using the existing designs and features selected from the database in steps 1 and
257 2. As a new design evolves, the PFLF is generated from the FLPF and the design type prior. The PFLF
258 can be updated in response to events to provide feature suggestions at interactive speeds.

259 4 Model Development

260 In this section, a model is mathematically developed for steps 3 to 6 from the process in Section 3. This
261 will allow for both predictions on feature type together with its spatial position.

262 4.1 Overview

263 We model the process of a designer choosing to place features in specific locations as a Marked Point
264 Process (MPP). As such we can view historical designs as a realization from this process. Consider a
265 design denoted by d_i , which comprises n_i features (not necessarily unique) and for each feature, which
266 we denote with m , we have an associated location described by its (x, y) coordinates. We express the
267 design as an unordered set of coordinates and features with $d_i = \{(x_1, y_1, m_1), \dots, (x_{n_i}, y_{n_i}, m_{n_i})\}$. We
268 restrict our designs to 2 dimensions expressed as (x, y) coordinates for simplicity but the method is easily
269 generalizable to higher dimensions.

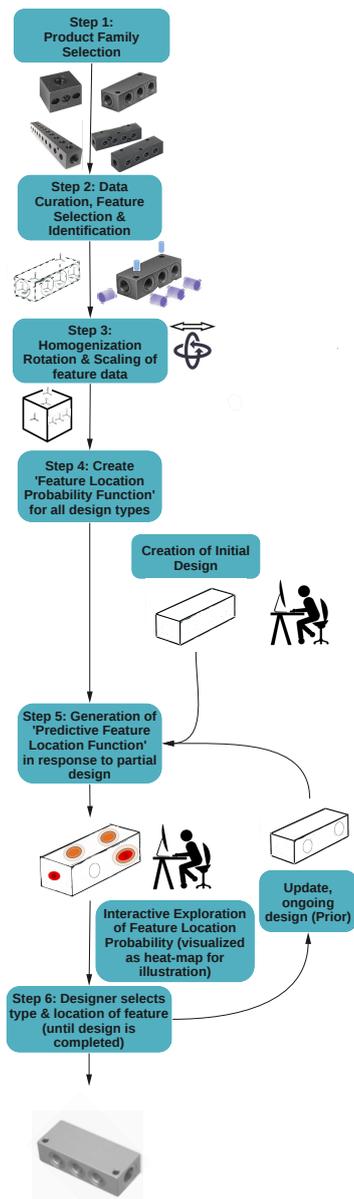


Figure 3: System schematic for feature prediction through Bayesian updating.

270 At the core of a MPP is the intensity function, which describes the probability of a feature being
 271 placed in a particular location. Let $\lambda(x, y, m | \underline{c})$ denote the intensity function of the process at location
 272 specified by the (x, y) coordinates, for feature m given the designer has already made choices of places
 273 various features at locations captured in the matrix \underline{c} . A characteristic of this function is that if we
 274 integrate the intensity over the whole (x, y) plane then we obtain the expected number of features m in
 275 the design. Moreover, we can express this as a probability density function, i.e. $f(x, y, m, | \underline{c})$ given in
 276 Eq. (1), to describe the next choice made by the designer by normalizing it so that it integrates to 1. This
 277 function can then be used to rank features based on their likelihood of being placed at specific locations
 278 to provide appropriate decision support to the designer.

$$f(x, y, m, | \underline{c}) = \frac{\lambda(x, y, m | \underline{c})}{\iint_{\forall x, y} \lambda(x, y, m | \underline{c}) dx dy} \quad (1)$$

279 Engineering designs possess dependency structures unlike other fields of MPP study so ‘off the shelf’
 280 models for intensity functions are not available. Dependency refers to the association of choices, such
 281 that placing one feature in a location increases or decreases the likelihood of other features in various
 282 locations. Poor choices of dependency models can result in uninformative inference at best and misleading
 283 inference at worst. In typical spatial or temporal point process applications, self-exciting models are used
 284 to capture local dependency where the realization of one point increases the likelihood of nearby points
 285 being discovered. In design, choosing a feature for a location can have ramifications for distant locations
 286 due to a need for symmetry for example. We develop a methodology for characterizing such dependency.

287 Many designs may be a collection of few choices, so while there may exist a large database of historical
 288 designs there are small sample sizes on which to infer the dependency structure. Inference is made more
 289 challenging with an extensive set of features from which to choose.

290 We propose a non-parametric approach to estimating the intensity functions that will provide a
 291 foundation on which to develop decision support, estimated from the data on historical designs. Kernel
 292 density estimates (KDE) consists of modeling the intensity function of a point process through assigning
 293 a kernel, e.g. the Normal distribution, centred at each location where a point has been realized, often
 294 resulting in a multi-modal probability model to describe the likelihood of discovering points. Typically
 295 the kernel density requires the analyst to choose a value for the smoothing parameter (e.g. in the case of
 296 the Normal kernel density this would correspond to the standard deviation for each density used).

297 In Section 4.2 we will develop the non-parametric model for the density function based on KDE from
 298 historical designs. In Sections 4.3 and 4.4 we will outline a Bayesian updating mechanism that will
 299 show how the density function changes as the designer makes further choices and as such so too will the
 300 decision support. In Section 4.5 we will derive metrics to characterize the dependency structure implied
 301 by these modeling assumptions. Finally in Section 4.6, we will consider transformation that we can make

302 on historical design data to improve predictions, specifically, re-scaling and rotating the data.

303 4.2 Model Description

304 We assume that a new design will be similar in some sense to historical designs but not necessarily
 305 identical. As such, prior to commencing an assessment should be made of the historical data that will
 306 be used to assess its suitability. Assuming we have a catalog of n historical designs that are appropriate
 307 for the decision support then we consider that there are n types of design and the current design under
 308 construction will belong to one of these types. We will estimate the density function for each type with
 309 the data available from each design. Following this we will apply a prior probability on the type of design
 310 being constructed based on the choices made.

311 Consider an historical design i for which there have been $n_{i,m}$ choices of feature m . Using a KDE
 312 approach to estimate the probability density function for design of type i with respect to feature m we
 313 have the density given in Eq. (2)

$$f(x, y|m, i) = \begin{cases} \frac{1}{c_{xy}} & n_{i,m} = 0 \\ \frac{\sum_{j=1}^{n_{i,m}} \phi_j(x, y; \mu_{x,j} = x_{i,m,j}, \mu_{y,j} = y_{i,m,j}, \sigma)}{n_{i,m}} & n_{i,m} \geq 1 \end{cases} \quad (2)$$

314 Where $\phi_j(\cdot)$ is a bivariate Normal density function, $\mu_{x,j}$ is the mean of the x variable in $\phi_j(\cdot)$, $\mu_{y,j}$
 315 is the mean of the y variable in $\phi_j(\cdot)$, $x_{i,m,j}$ is the location on the x coordinate in design i of the j^{th}
 316 occurrence of feature m , $y_{i,m,j}$ is the location on the y coordinate in design i of the j^{th} occurrence of
 317 feature m , σ is the standard deviation for both x and y , although one could assume a more elaborate
 318 covariance structure if appropriate, and c_{xy} is a normalizing constant to ensure the density integrates to
 319 1. It is worth noting that one could substitute other kernel density functions in if more appropriate, we
 320 only require it to possess all the characteristics of a bivariate probability density function.

321 We have assigned a uniform distribution over the plane for situations where that feature has not
 322 appeared in design i . It may be desirable to remove this, if one did not want to permit certain features
 323 for particular design types.

324 Essentially, the resulting density is a collection of Normal densities centred about observed locations
 325 and the standard deviation parameter controls for the allowable variation from the historical design to
 326 be considered similar.

327 Let I be the random variable describing the design type that the designer is developing and M to be
 328 the random variable describing the next feature to be chosen. To express the unconditional probability
 329 density function we first define three indicator functions to denote design type, feature, and presence in
 330 Eq. (3).

$$\delta_i = \begin{cases} 1, & I = i \\ 0, & I \neq i, \end{cases} \quad \delta_m = \begin{cases} 1, & M = m \\ 0, & M \neq m, \end{cases} \quad \delta_{n_{im}} = \begin{cases} 1, & n_{im} \geq 1 \\ 0, & n_{im} = 0 \end{cases} \quad (3)$$

331 We denote the probability of a feature m appearing in design type i with $p_{i,m}$ and the probability of
 332 the design being of type i with $\pi(i)$. Combining these, the full probability density function describing
 333 the likelihood of a feature m being located at (x, y) and the design being of type i is given in Eq. (4).

$$f(x, y, M = m, I = i) = \sum_{i=1}^{i_{max}} \sum_{m=1}^{m_{max}} \delta_i \pi(i) \delta_m p_{i,m} \left(\delta_{n_{im}} \frac{\sum_{j=1}^{n_{i,m}} \phi_j(x, y; \mu_{x,j} = x_{i,m,j}, \mu_{y,j} = y_{i,m,j}, \sigma)}{n_{i,m}} + (1 - \delta_{n_{im}}) \right) \quad (4)$$

334 The design type is a latent variable used to capture the dependency between the features and locations.
 335 By summing the density function across all possible values of I we obtain the distribution for location
 336 and feature only, given in Eq. (5).

$$f(x, y, M = m) = \sum_{i=1}^{i_{max}} \sum_{m=1}^{m_{max}} \pi(i) \delta_m p_{i,m} \left(\delta_{n_{im}} \frac{\sum_{j=1}^{n_{i,m}} \phi_j(x, y; \mu_{x,j} = x_{i,m,j}, \mu_{y,j} = y_{i,m,j}, \sigma)}{n_{i,m}} + (1 - \delta_{n_{im}}) \right) \quad (5)$$

337 Similarly, we express each marginal distribution in Eqs. (6–9).

$$f(x) = \sum_{i=1}^{i_{max}} \sum_{m=1}^{m_{max}} \pi(i) p_{i,m} \left(\delta_{n_{im}} \frac{\sum_{j=1}^{n_{i,m}} g_j(x; \mu_{x,j} = x_{i,m,j}, \sigma)}{n_{i,m}} + (1 - \delta_{n_{im}}) \right) \quad (6)$$

$$f(y) = \sum_{i=1}^{i_{max}} \sum_{m=1}^{m_{max}} \pi(i) p_{i,m} \left(\delta_{n_{im}} \frac{\sum_{j=1}^{n_{i,m}} g_j(y; \mu_{y,j} = y_{i,m,j}, \sigma)}{n_{i,m}} + (1 - \delta_{n_{im}}) \right) \quad (7)$$

$$f(M = m) = \sum_{i=1}^{i_{max}} \pi(i) p_{i,m} \quad (8)$$

$$f(I = i) = \pi(i) \quad (9)$$

338 Where $g_j(\cdot)$ is a univariate Normal density function.

339 4.3 Probability of a Feature Being Selected for a Design

340 Given the total number of incidences within a design we assume the number of incidences of each possible
341 feature for a design is a realization from a multinomial distribution. Moreover, we assume the underlying
342 probabilities associated with each feature vary across design types. Under such a modeling assumption
343 a natural estimator of the probability of a feature being selected for a design of a particular type would
344 be the observed frequency on similar designs from the class. However, given that we have at most one
345 design for each type we are likely to produce poor inference due to small samples. Moreover, we are likely
346 to be faced with a large number of features with zero events data resulting in an estimated probability
347 of 0. This creates a particular issue for the decision support being developed, as all historical designs
348 that did not possess all the features chosen for a current design would be ruled out as candidate design
349 types through Bayesian updating. As such, allowing for non-zero probability estimates would permit the
350 inclusion of candidate design types even if they do not include all the features chosen at some point in
351 the design process. For a discussion on alternative estimation methods for zero event data, see [32].

352 We propose using an uninformative prior distribution, where a uniform prior distribution is assumed
353 on each probability and subsequently updated on the data. As the probabilities must sum to 1, the
354 uniform assumption implies a Dirichlet prior distribution. This is a common pairing with the multinomial
355 distribution as it provides a flexible distribution that is convenient to use computationally. This results
356 in the following estimate for the probability, given in Eq. (10).

$$p_{i,z} = w \frac{\beta_z}{\beta} + (1 - w) \frac{k_{i,z}}{k_i} \quad (10)$$

357 Where $w = \beta/(\beta+k_i)$, $k_{i,z}$ gives the number of features in design i of type z , $k_i = \sum_{\forall z} k_{i,z}$, $\beta_z =$
 358 $\sum_{\forall i} k_{i,z}$, and $\beta = \sum_{\forall z} \beta_z$.

359 We see that $p_{i,z}$ is a weighted average of the observed frequency $k_{i,z}/k_i$ and the prior mean. The
 360 weight applied to the frequency increases as the number of features chosen for design i increases, i.e. k_i .

361 4.4 Bayesian Updating

362 Every choice made by the designer provides information concerning the type of design being constructed,
 363 i.e. to which historical design is it similar. We will model this learning through Bayesian updating. As
 364 described in Section 4.2 we have a probability distribution, i.e. $\pi(i)$, which describes the uncertainty
 365 concerning the design type. In this section we present a Bayesian updating of this distribution based
 366 on design choices. Assume that the designer has made n_k choices then the posterior distribution for the
 367 design type is updated as in Eq. (11).

$$\pi(i|\mathcal{C}) = \frac{\prod_{k=1}^{n_k} f(x_k, y_k, M = m_k, I = i)}{\sum_{i=1}^{i_{max}} \prod_{k=1}^{n_k} f(x_k, y_k, M = m_k, I = i)} \quad (11)$$

368 This posterior is then used in the predictive distribution, given in Eq. (12).

$$f(x, y, M = m|\mathcal{C}) = \sum_{i=1}^{i_{max}} \sum_{m=1}^{m_{max}} \pi(i|\mathcal{C}) \delta_m p_{i,m} \left(\delta_{n_{im}} \frac{\sum_{j=1}^{n_{i,m}} \phi_j(x, y; \mu_{x,j} = x_{i,m,j}, \mu_{y,j} = y_{i,m,j}, \sigma)}{n_{i,m}} + (1 - \delta_{n_{im}}) \right) \quad (12)$$

369 This function can be used to provide inference on the relative likelihood of features being located on
 370 specified positions through comparing ratios.

371 4.5 Dependency Structure

372 The moments of the model are easily obtained through conditional expectation arguments resulting in the
 373 expectations given in the Supplemental Material S1. Through setting $\delta_{n_{im}} = 1$ for all designs we would
 374 obtain the moments anticipated in the historical designs, however, for our model we have accommodated
 375 the possibility of features appearing in design types which are not present in the associated historical
 376 design.

377 The moments can then be used to construct measures such as correlation between the (x, y) coor-

378 dinates. However, while commonly used such measures are limited within our context as they focus on
379 the linear relationship between only two variables. We may wish to consider more general settings such
380 as non-linear relationships as well as 3D designs or even the dependency between features and locations.
381 For this we use the mutual information (MI) measure, which we denote by ω , to assess dependency. The
382 concept of mutual information is linked to the entropy of a random variable, which quantifies the expected
383 amount of information held in a random variable. The mutual information measure is considering the
384 information gain from modeling the joint distribution rather than assuming each variable is independent.

$$\omega = \mathbb{E} \left[\ln \left(\frac{f(X, Y, M)}{f(X)f(Y)f(M)} \right) \right] \quad (13)$$

385 This can be re-expressed as in Eq. (14), which is simply the entropy of the joint distribution minus
386 the sum of the entropy for all the marginals.

$$\omega = \mathbb{E} \left[\ln (f(X, Y, M)) \right] - \mathbb{E} \left[\ln (f(X)) \right] - \mathbb{E} \left[\ln (f(Y)) \right] - \mathbb{E} \left[\ln (f(M)) \right] \quad (14)$$

387 Clearly, $\omega = 0$ if $f(x, y, m) = f(x)f(y)f(m)$, i.e. if the variables are independent. Moreover, it can be
388 shown that as dependency increases so too does the measure. This measure can be useful for comparing
389 dependency between various subsets of designs, noting that the stronger the dependency the better the
390 predictions will be. Some analysts prefer to transform this measure to bound it within $(0, 1)$ and as such
391 use the transform $\dot{\omega} = \sqrt{1 - e^{-2\omega}}$ [33]. The joint and marginal distributions required to calculate the MI
392 are provided in S2.

393 4.6 Re-Scaling and Rotating

394 Generally, the higher the degree of homogeneity in the comparator pool of data, then the greater the
395 accuracy in the prediction [30] and as such pre-processing the relevant historical data to achieve greater
396 homogeneity may be desirable. We consider re-scaling and rotating the data for each as a means to
397 achieve this. However, such transformations may not always be beneficial as key information may be
398 lost that helps identify the most similar historical designs. The advantage of such transformations are
399 through identifying regions where specific features are highly likely to be located for a large number of
400 design types. The disadvantage can be blurring distinctive characteristics between design types and as
401 such it will take longer for the process to learn precisely to which design type it belongs.

402 Re-scaling can be achieved through stretching or compressing a design to the unit cube, so that the
403 length of each dimension is re-scaled such that 0 is the minimum and 1 the maximum in that dimension for
404 that design. If such a transformation is performed, care must be taken in interpreting distance between
405 two points as the scales would not be the same between dimensions. Rotations can be the result of
406 a non-standardized axis used with designs. As such, through rotation the data we are constructing a

407 common axis, which may reveal more similarity across the designs. Rotation data is achieved through
 408 matrix multiplication of the data set. For example, in a 2 dimensional design, every rotation around the
 409 origin in a counter clockwise direction can be represented with the matrix R shown in Eq. (15).

$$R = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix} \quad (15)$$

410 When the data are multiplied by R we obtain the new coordinates as in Eqs. (16).

$$\begin{aligned} x_j(\theta) &= \cos(\theta)x_j - \sin(\theta)y_j \\ y_j(\theta) &= \sin(\theta)x_j + \cos(\theta)y_j \end{aligned} \quad (16)$$

411 An analyst could decide upon rotation and rescaling based on visual inspection. However, for a
 412 more rigorous approach we would need to measure the distance between designs and seek to minimize
 413 it. Following the approach proposed by [31] we use the Kullback-Leibler (KL) divergence measure to
 414 assess the difference between designs. Using the superimposition of all designs as an average design we
 415 can measure the difference of each design to the average and seek to minimize it.

416 The KL divergence measure of design type v to u , denoted by $D_{KL}(D_u \parallel D_v)$ is given in Eq. (17).

$$\begin{aligned} D_{KL}(D_u \parallel D_v) = & \sum_{m=1}^{m_{max}} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} P_{u,m} \left(\frac{\sum_{j=1}^{n_{u,m}} \phi_j(x, y; \mu_{u,x,j} = x_{u,m,j}, \mu_{u,y,j} = y_{u,m,j}, \sigma)}{n_{u,m}} \right) \\ & \ln \left(\frac{P_{u,m} \left(\frac{\sum_{j=1}^{n_{u,m}} \phi_j(x, y; \mu_{u,x,j} = x_{u,m,j}, \mu_{u,y,j} = y_{u,m,j}, \sigma)}{n_{u,m}} \right)}{P_{v,m} \left(\frac{\sum_{j=1}^{n_{v,m}} \phi_j(x, y; \mu_{v,x,j} = x_{v,m,j}, \mu_{v,y,j} = y_{v,m,j}, \sigma)}{n_{v,m}} \right)} \right) dx dy \quad (17) \end{aligned}$$

417 This is re-expressed in Eq. (18).

$$D_{KL}(D_u \parallel D_v) = \mathbb{E} \left[\ln \left(p_{u,m} \left(\frac{\sum_{j=1}^{n_{u,m}} \phi_j(x, y; \mu_{u,m,j} = x_{u,m,j}, \mu_{u,y,j} = y_{u,m,j}, \sigma)}{n_{u,m}} \right) \right) \right] - \mathbb{E} \left[\ln \left(p_{v,m} \left(\frac{\sum_{j=1}^{n_{v,m}} \phi_j(x, y; \mu_{v,m,j} = x_{v,m,j}, \mu_{v,y,j} = y_{v,m,j}, \sigma)}{n_{v,m}} \right) \right) \right] \quad (18)$$

418 Where the expectation is taken with respect to the distribution with u . Expressing this as an expect-
419 tation provides a computational advantage, as a closed form analytical solution is not available, we can
420 conduct Monte Carlo simulations with the distribution of u and evaluate the average of the expression.
421 In sum, we can transform each design through rotation and re-scaling to minimize the KL divergence of
422 the mean design to the design in question.

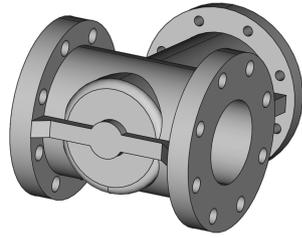
423 4.7 Summary

424 Section 4 has outlined the underlying model and process to support the prediction of features given
425 location. This can be used with an interactive CAD system, where the cursor sits in a location described
426 by its coordinates and the recommended feature is suggested. In Section 5 we apply this to a data set.

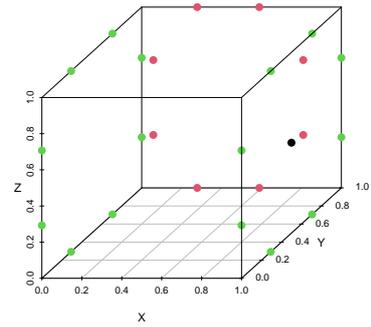
427 5 Case Study

428 To allow an intuitive, visual understanding of the proposed process we have chosen to use a set of 513
429 mechanical valve designs. The structure of the valve bodies have obvious regularities with circles around
430 the valve's flanges together with other functional holes. An unordered set of hole diameters and associated
431 (x, y, z) coordinates were extracted for each valve body from the B-rep of the CAD design using the Twig
432 match algorithm [34]. Further details are provided in [31]. An example valve design is shown in Fig. 4a
433 with the extracted hole features, scaled to $[0, 1]$, shown in Fig. 4b.

434 In this analysis, the aim is to predict the sequential addition of hole features and their position
435 given the state of the current design, with the focus on features occurring on the same surface plane i.e.
436 predicting a hole diameter on the flange surface.



(a)



(b)

Figure 4: Example of a valve body from CAD design database. Figure 4a shows an image of a valve design and Fig. 4b the scaled positions of the extracted hole features. The different colors are used to identify the different diameters of the holes

437 5.1 Scaling

438 To facilitate prediction the feature coordinates of each design were scaled to the unit cube – each dimension
 439 was scaled to $[0, 1]$ – and additional rescaling was required on each cube surface so that the features
 440 retained their geometric shape. Models were then estimated using the features and feature positions
 441 which were positioned on the surface of the cube, one surface at a time. For example, after scaling, the
 442 data were subset to analyze the features on the $x = 0$ face. Figure 5 shows the superimposition of the
 443 scaled feature coordinates from all designs on the $x = 0$, $y = 0$ and $z = 0$ faces (some jitter of the points
 444 has been added to the Figure to aid feature discrimination).

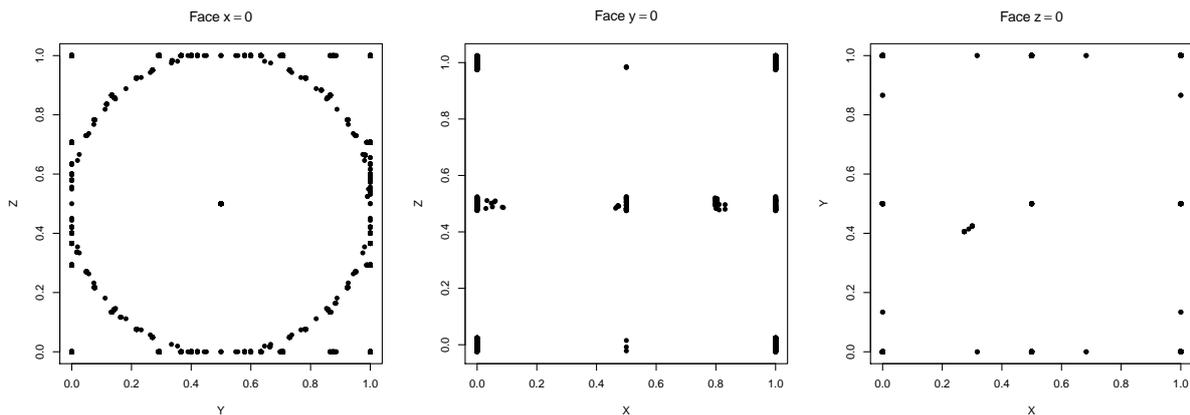


Figure 5: Superimposition of scaled hole feature coordinates for all the data on the specified face of the cube.

445 5.2 Kernel Density Estimation

446 The KDE were estimated across each face of the cube. This was done by dividing the face into N by N
447 regular grid positions and then estimating the kernel density at each grid position for each feature in all
448 designs. A Normal kernel with user specified standard deviation was applied, in this analysis chosen to
449 be 0.05, and the density on each dimension calculated independently. If a design did not have a specific
450 feature that was present in the database pool of features then a Uniform probability across the grid
451 of positions was assumed. This allows for predictions to be generated on a new design which is using
452 a combination of hole features that have not been previously observed. The KDE outputs a density
453 estimation at each position in the N by N grid for every feature in the database.

454 5.3 Evaluation

455 Both the correctness of the FLPF and the predictive accuracy of the PFLF were assessed using 10-
456 fold cross-validation. The kernel density across the features was estimated using the training data and
457 then evaluated on the designs in the test set. Each test design, which contain hole feature labels and
458 their coordinates, was evaluated one at a time. Three measures of predictive performance were used;
459 the distance from the observed feature coordinate to the nearest predicted mode was calculated using
460 two approaches, and reciprocal rank was used to evaluate how accurately a feature was predicted at its
461 observed position. Further details are provided in the Supplemental Material S3.3.

462 An illustration of the Bayesian updating and predictions on one test design is provided in the Sup-
463 plemental Material S3.

464 5.4 Results

465 Figure 6 illustrates the aggregated results from the cross-validation. The x-axis indicates how many holes
466 have been added to a new design (e.g. if a test design has four features, then the predictive densities,
467 distances to mode and ranks are calculated after sequentially adding 0, 1, 2, or 3 holes to the new design).
468 The y-axis gives the mean distance to the mode, either on the raw scale or in grid steps or the mean
469 reciprocal rank. The red triangle gives the mean across the ten folds. The performance of the predicted
470 rank of suggestions is shown in the third figure; the range of values is from zero – poor suggestions, to
471 one – perfect suggestions.

472 As expected an initial improvement is observed in the distances to the nearest predictive mode as
473 additional features were added to each new test design, however, there is a clear pattern of extreme values
474 within all figures which results in a decrease in performance as additional holes are added to a design.
475 This can be explained, as within each test fold there are a few designs which are unlike anything in the
476 training set and thus the KDE does not provide reasonable predictions.

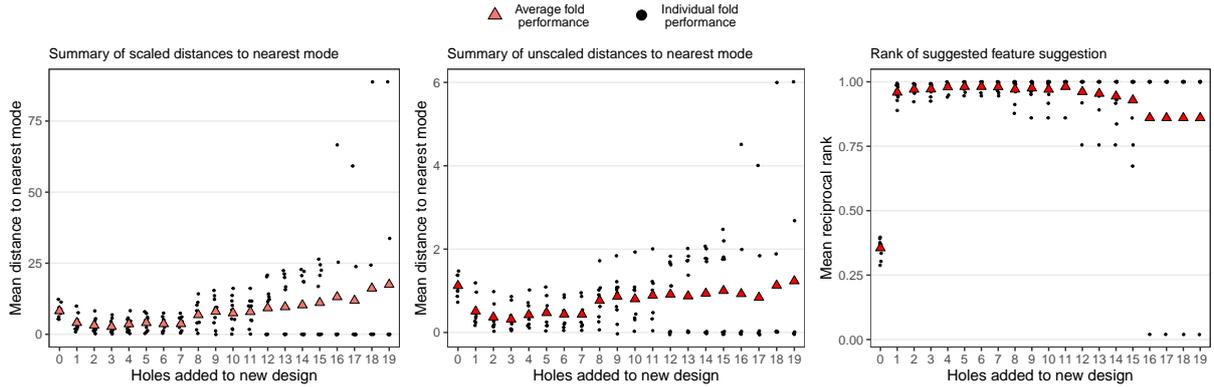


Figure 6: Cross-validation predictive performance. The x-axis indicates how many holes have been added to a new design and the y-axis records the average distance or rank. The black points give the within fold average measure and the red triangle the average performance across the folds.

477 Some examples follow. All the designs in which a specific feature occurs in the training dataset may
 478 have a different number of feature instances than observed in the test design. In one example, within
 479 a training dataset instance, designs with the hole diameter “33.0” have between 3 and 13 instances on
 480 the flange plane, however it occurs 16 times within a test design. This results in all the predictive
 481 modes being slightly offset, as seen in the Supplemental Material S6.1. This modeling framework cannot
 482 infer the coordinates for features even though they are still placed within the same circular orientation.
 483 Another example is that all of the designs with a specific feature in the training dataset are positioned
 484 differently than those in the test design. For an example shown in the Supplemental Material S6.2, the
 485 hole diameter “35.0” was used as a central bore hole in the training data designs, however, it was used as
 486 the bolt connector within the new test design. Therefore as additional holes are added to the new design
 487 the updated predictive density provides little information. Clearly the order in which the features are
 488 added to a new design will affect the predictive density, particularly when there are multiple types (hole
 489 dimensions) of feature, and this can impact the quality of predictive guidance. For an example shown in
 490 the Supplemental Material S6.3, there are 30 designs within a training dataset instance that contain the
 491 “22.0” diameter, but only one of these also has the additional “17.29” diameter. The early selection of
 492 the “17.29” feature adds more probability weight onto the single design in the training set, and it takes
 493 several further additions for the predictions to improve.

494 The predictive performance of the method was re-evaluated omitting the 24 unusual designs from the
 495 test datasets and the results are shown in Fig. 7. This is done to examine the predictive performance
 496 of the model for a designer who remains within the catalog of previous designs. It can be seen that the
 497 predictive performance improves as more features are added to a new design. This again indicates that
 498 the utility of the method is dependent on the designs forming a homogeneous set. The folds with larger
 499 values can be explained by the ordering of the features entering in to the design, as illustrated by the
 500 example given in the Supplemental Material S6.3. While the distance to the nearest predictive mode

501 may be small, there remain extreme values in the rank predictions. This indicates that while a feature
 502 is expected at a position, our model has been unable to predict the specific feature, and so suggests that
 503 the feature added at this position is unusual given those observed at in the training data.

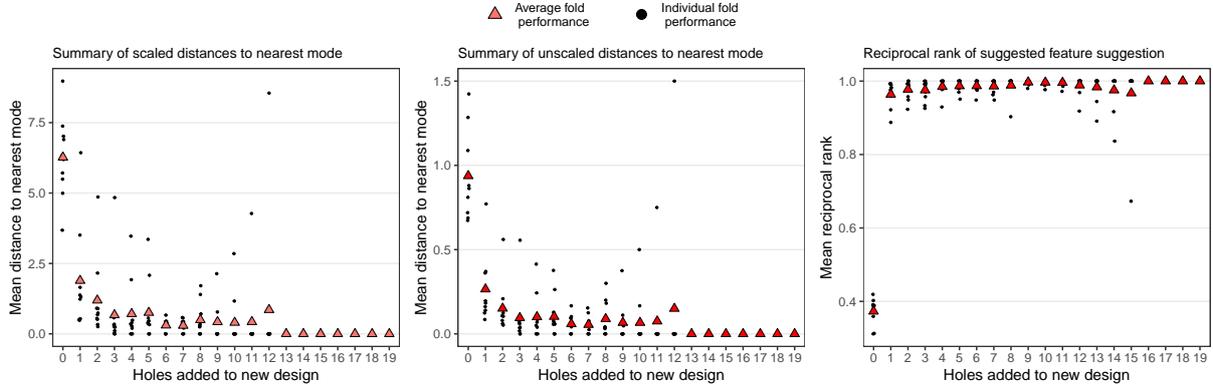


Figure 7: Cross-validation predictive performance after omitting unusual designs from the test set. The x-axis indicates how many holes have been added to a new design and the y-axis records the average distance or rank. The black points give the within fold average measure and the red triangle the average performance across the folds.

504 The prediction results for the features on the $y = 0$ and $z = 0$ cube face are provided in the Supple-
 505 mental Material S7. Performance is similar to the $x = 0$ face with predictions improving as additional
 506 features are added to the design. Within the features on these faces, there are two designs that are unlike
 507 any of the other designs; the effect of this is more apparent on the $y = 0$ face. Again removal of the
 508 outlier design resulted in improved statistics (results not provided).

509 5.5 Association Measure

510 The utility of the method is supported by the homogeneity of the design database. Section 4.5 described
 511 how the mutual information could be used to provide some measure of the expected dependence in a
 512 database, however, there is no analytic solution to Eq. (14) for our non-parametric model. We therefore
 513 estimate this measure through a simulation exercise. We denote i_{max} as the number of designs in the
 514 database, m_{max} as the number of unique marks (features) and $n_{i,m}$ as the number of marks of type m
 515 in design i . The probability of randomly choosing design i is given by $q_i = 1/i_{max}$ and the probability of
 516 selecting mark m given that design i was selected is defined by $p_{i,m} = n_{i,m} / \sum_{\forall m} n_{i,m}$.

517 For a given design database, a representative random sample of feature instances is generated, using
 518 the following steps,

- 519 1. Uniformly sample a design i from the set of designs in the database with probability q_i
- 520 2. Randomly sample a feature type m from design i with probability $p_{i,m}$
- 521 3. Sample a single instance of feature type m , as there may be multiple instances of feature m within
 522 design i .

- 523 4. Take a random sample from the Normal kernel with mean at the (x, y) coordinates
- 524 5. Repeat many times

525 An estimate of the MI can then be calculated using the expressions given in the Supplemental Material
526 S2 where the KDE of the designs in the database are evaluated at the sampled coordinates.

527 The dependence structure within our sample database was estimated using this method across the
528 ten training cross-validation datasets, and the resulting MI scores had mean 1.43 and standard deviation
529 0.03. This equates to a scaled $\hat{\omega} = 0.97$ indicating that there is strong dependence within the data and
530 thus we would expect predictions to be good. For comparison, a null distribution was estimated for
531 the statistic on the same training data, permuting the feature instance and generating the coordinates
532 randomly from the Uniform $[0,1]$ distribution, so that features were no longer aligned with specific designs
533 or coordinates. This gave a MI of 0.62 (0.03), and scaled value of 0.84. A second smaller simulation of
534 randomly generated designs and feature coordinates revealed that smaller samples produced higher MI.
535 As sample size increased then the MI decreased to zero, the theoretical value for independence. It
536 would therefore be useful for practitioners to evaluate the MI on randomly permuted data to support
537 interpretation of the MI score on the design database.

538 5.6 Rotation

539 The more similar the designs in the database the stronger the signal for making predictions. However,
540 different designs may have been created with a different orientation. Section 4.6 described how the
541 KL measure could be used to rotate one design to minimize the probabilistic differences between them.
542 There is not an analytic solution to Eq. (18) for our non-parametric model but we can minimize the KL
543 divergence between two designs u and v , $D_{KL}(D_u || D_v)$, by finding the angles of rotation that maximize
544 the second term by a simulation design embedded in an optimization routine. As the KDE may be a
545 noisy function, a global optimization routine should be used, although a brute force search is feasible in
546 2D.

547 The rotation can be implemented as follows. The design database contains feature instances that
548 are assumed to be representative of the underlying orientation of the designs and we consider this as a
549 single average design. Each design would then be orientated in turn, to this average design, excluding the
550 design getting rotated from the pool. First multiple random draws are simulated from the Normal kernel,
551 with means equal to the feature positions of the design to be rotated. Then the KL measure is calculated
552 between the samples and the average design. The new design is then rotated and re-sampled until the KL
553 measure is minimized. An illustration of the 2D rotation of a part is given in the Supplemental Material
554 S5.

5.7 Choice of Standard Deviation

The choice of standard deviation (bandwidth) of the Normal kernel determines how spread out the predicted density is around the training data observations. Figure 8 shows the predicted density for the same training data under different kernel standard deviations. There are two main approaches one could take in determining a suitable value for this parameter. It may be that prediction is the only desirable performance measure and then through cross validation exercises an optimal value can be identified. Alternatively, with more emphasis on the prospective nature of this decision support, to facilitate the determination of designs that are similar to historical designs, this parameter can be used as a controlling lever, whereby small values will result in predictions that are very close to previous designs.

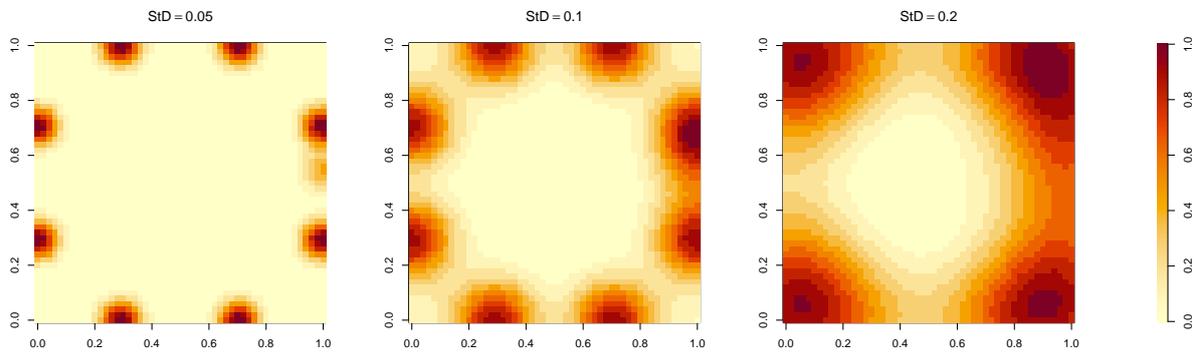


Figure 8: Predictive density calculated under a Normal kernel using different standard deviations. The colourbar legend describes how the image colour maps to the data, with dark red indicating regions of higher predictive density

5.8 Discussion of Results

We have described how to implement the proposed process for predicting the type and location of the features that might be added during an engineering design process. We evaluated the method on a data set of real designs through a cross validation process. In 90% of the evaluation runs the feature's actual location and the prediction (once at least one feature had been selected) were very close (i.e. within 0.5 grid space on average – 1% of the normalized range of the part). When more features were added to the design the accuracy of the predictions improved. This observation can be clearly seen in the ranking of the predicted features (i.e. an ordered list of the most to least likely features to occur at a given location). If four features had been selected (i.e. added to the design) the subsequent features selected were, on average, ranked in the first 25% of the list of suggestions. This increased to the top 10% once eight features had been selected.

This behavior reflects the nature of the commercial product families which formed the dataset. These have frequently repeated sets of features at standardized positions within a design, and so after one choice has been made then subsequent choices can be predicted with a high probability. In other words portfolios

578 of mechanical designs have strong dependence in data that results in strong predictive performance.

579 Figure 9 shows two components with a feature prediction “heat map” manually superimposed on
580 to their faces. These hotspots indicate where features (regardless of their type and size) are frequently
581 located in the training data. For each figure, the left-hand side plot presents the complete part, the
582 central plot presents the predictive density from the training pool when no features have been added to
583 the part, and the right-hand side plot the updated density given that the new features have been added
584 to the design.

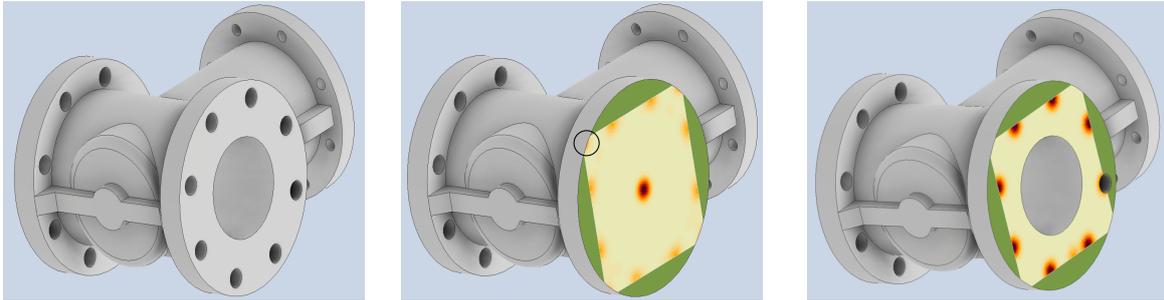
585 Figure 9b presents an example of how the predictive density changes once a feature has been added to
586 the design. Before a choice has been made, the greatest predictive density is placed on the four corners
587 of the feasible box, the normalized range of the scaled prediction region, indicating that this orientation
588 is most common in the training dataset. Once a feature of specific size and position has been added, the
589 predictive density changes to favor a pattern of six holes away from the corners. Red and blue circles
590 are used in Fig. 9b to highlight how the predictive density in these regions change given some feature
591 addition. Although such heat-maps give an intuitive overview the prediction results can be presented to
592 a user in several different ways. For example, given a location (e.g. the user’s cursor) a list of feature
593 suggestions (ranked in order of their likelihood) could be generated.

594 The heat maps also illustrate the need for further research into user interfaces that allow the designer
595 to control the choice of training data (used to generate the predictions) and the scaling/mapping of
596 the results onto new designs. In the case study, the feature coordinates were normalized to boundaries
597 determined by the extent of feature locations within the training dataset. For example, in Fig. 9a
598 predictions were generated across the unit cube and so there is non-zero density in the corners (highlighted
599 by black circle in central panel of Fig. 9a), whereas if a unit circle had been used to normalize the feature
600 locations the result would be more appropriate to the shape by omitting the truncated predicted region
601 corners. An obvious artifact of the current approach to normalization is that there will be regions on
602 a face (colored green in Fig. 9) that are beyond the geometric extent of the features used to train the
603 prediction system. Due to the restrictions that we have imposed through this mapping, predictions were
604 not generated outside the range of the normalized region. However, depending on the choice of kernel,
605 one could extrapolate beyond these boundaries for a new design, but as with any extrapolation, these
606 require stronger assumptions.

607 This could be mitigated by filtering suggestions that are physical or functionally feasible before pre-
608 senting them as options to the designer. The development of effective filters would also enable the
609 geometric limits on feature prediction to be determined in a manner most appropriate to the application
610 (e.g. a part bounding box, or specific planes). However, work is required on generic scaling functions to
611 support this, for example the top flange of the design-part provided in Fig. 9b is rectangular and so the
612 predictive density requires to be mapped back to the original part dimension from the unit cube used to

613 generate predictions.

(a) Part-design A



(b) Part-design B

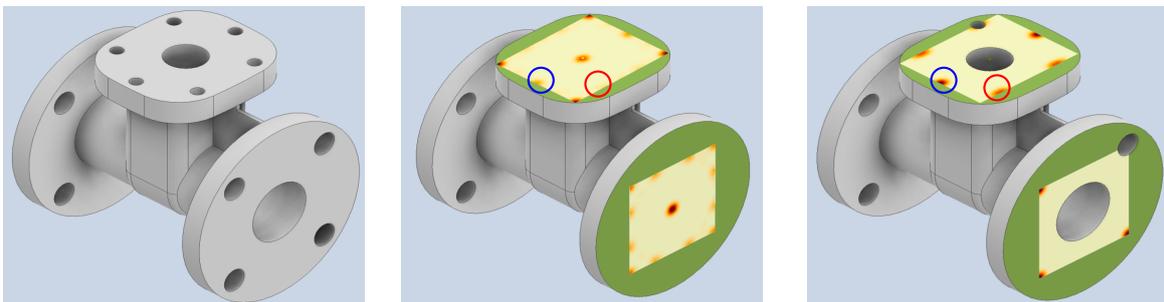


Figure 9: The feature predictions were mapped onto two incomplete designs. Regions of higher predictive density are illustrated by dark red. From left to right, the plots show the complete part, the predictive distribution when no features have been added, and so is completely informed by the training data, and lastly the predictive distribution after the addition of a central bore hole and one bolt hole feature. The updated distribution displays the increased level of belief in the positions of subsequent features, with areas of higher density and lower variance. The black circle in (a) is used to indicate the prediction region which has been truncated and the red and blue circles in (b) are used to highlight the changes in the predictive density once a choice has been made. The green sections represent areas outside the normalized prediction region.

614 6 Summary and Conclusions

615 The aim of this research was to “Define a computational framework that can support an interactive
616 design process with suggestions of features based on three inputs: a knowledge of existing designs; the
617 state of an emerging design and a location on the surface of the emerging design.” The authors believe
618 that system described meets this goal and has established how the feature content of mechanical designs
619 can be amalgamated and transformed into a likelihood function that defines the probability of particular
620 design features occurring at specific locations on a model.

621 The work has not only demonstrated that the architecture of the proposed system is viable but also
622 established that the computations can be done quickly enough to support a dynamic design process. For
623 example, the prototype system can respond to a given mouse location at interactive speeds (i.e. ms)
624 and consequently could support user interface functionality such as pop-up menus (customized to reflect

likely feature types and parameter values) or even ghosting images of possible features onto the cursor location as it moves to particular locations. In this way the engineer is free to ignore these selections in the same way a user of a predictive text system is able to adopt or dismiss suggestions when composing SMS texts.

The case-study illustrated the method using hole features, however, the feature set could be extended to include those with a more complicated geometry. Provided that a feature can be defined geometrically and hence extracted from the CAD design, the prediction method can be applied by considering such features as another type of mark. This would allow for modeling the dependence both between and across feature types.

6.1 Limitations

Like other predictive systems and there are inevitable limitations. Currently the system can only predict the likelihood of features occurring within the volume defined by the maximum extent of the features extracted from the training dataset. Understanding how these results can be generalized to support predictions across variable volumes, as well as optimal scaling of the normalized prediction region, is an area of further research. Additionally, while the method of data homogenization appears to be viable for product families with very regular structures (e.g. industrial valves or manifold blocks), its behavior with product families with more variable forms is not clear.

However, one of the features of all interactive predictive systems (that makes them viable) is that the user is always free to ignore suggestions that are wrong or out of context. In other words predictive systems do not have to provide perfect predictions all the time to be useful.

6.2 Future Work

Having established the fundamentals of the theory the authors intend to broaden the application to other datasets of mechanical component designs. This will allow the investigation of the methodology's ability to support multiple feature types and more geometrically varied product families i.e. the scaling of the normalized prediction region. The merits and implications of estimating the normalized prediction region using different kernels which can account for boundary effects will be studied. Considering MPP's beyond simple Euclidean geometry provide opportunities. The current focus of the project has been on providing decision support to a single engineer, and how such a system will support concurrent designs carried out simultaneously by distributed teams is a topic that requires further investigation.

Although this work has established the theoretical and computational foundations for a predictive system its utility will ultimately depend on how its user interface behaves. Although beyond the scope of this work follow-on projects will seek to incorporate the predictive functionality described in a commercial CAD system (via their API) and so allow a systematic assessment of the impact of predictive CAD on

658 design productivity to be undertaken.

659 **6.3 Acknowledgements**

660 This work was supported by the Engineering and Physical Sciences Research Council, UK [grant number
661 EP/R004226/1]. The dataset of hole features extracted from the valves models is available at [https://
662 doi.org/10.15129/c3ae80a3-b34e-43cc-b6e6-629eb5cab922](https://doi.org/10.15129/c3ae80a3-b34e-43cc-b6e6-629eb5cab922) (the URL will become active on journal
663 acceptance).

References

- [1] John E Ettlíe and Matthew Kubarek. Design reuse in manufacturing and services. *Journal of Product Innovation Management*, 25(5):457–472, 2008.
- [2] Rob Bracewell, Ken Wallace, Michael Moss, and David Knott. Capturing design rationale. *Computer-Aided Design*, 41(3):173–186, 2009.
- [3] Christopher McComb, Jonathan Cagan, and Kenneth Kotovsky. Mining process heuristics from designer action data via hidden markov models. *Journal of Mechanical Design*, 139(11), 2017.
- [4] Ayush Raina, Christopher McComb, and Jonathan Cagan. Learning to design from humans: Imitating human designers through deep learning. *Journal of Mechanical Design*, 141(11), 2019.
- [5] Suyu Hou and Karthik Ramani. Dynamic query interface for 3d shape search. In *International Design Engineering Technical Conferences and Computers and Information in Engineering Conference*, volume 46970, pages 347–355, 2004.
- [6] Silvia Schacht and Alexander Mädche. How to prevent reinventing the wheel?—design principles for project knowledge management systems. In *International Conference on Design Science Research in Information Systems*, pages 1–17. Springer, 2013.
- [7] John Ashburner and Karl J Friston. Voxel-based morphometry—the methods. *Neuroimage*, 11(6):805–821, 2000.
- [8] A. Getis, J. Lacambra, and H. Zoller, editors. *Spatial econometrics and spatial statistics*. Citeseer, 2004.
- [9] Siddhartha Chaudhuri and Vladlen Koltun. Data-driven suggestions for creativity support in 3D modeling. *ACM Transactions on Graphics*, 29(6), 2010.
- [10] Siddhartha Chaudhuri, Evangelos Kalogerakis, Leonidas Guibas, and Vladlen Koltun. Probabilistic reasoning for assembly-based 3D modeling. *ACM Transactions on Graphics (Proc. SIGGRAPH)*, 30(4), 2011.
- [11] Evangelos Kalogerakis, Siddhartha Chaudhuri, Daphne Koller, and Vladlen Koltun. A probabilistic model for component-based shape synthesis. *ACM Transactions on Graphics (TOG)*, 31(4):1–11, 2012.
- [12] Prakhar Jaiswal, Jinmiao Huang, and Rahul Rai. Assembly-based conceptual 3D modeling with unlabeled components using probabilistic factor graph. *Computer-Aided Design*, 74:45–54, 2016.
- [13] Lauren Lam, Sharon Lin, and Pat Hanrahan. Using text n-grams for model suggestions in 3D scenes. In *SIGGRAPH Asia 2012 Technical Briefs*, pages 1–4. 2012.

- 695 [14] Matthew Fisher, Manolis Savva, and Pat Hanrahan. Characterizing structural relationships in scenes
696 using graph kernels. *ACM Transactions on Graphics*, 30(4), 2011.
- 697 [15] Siddhartha Chaudhuri, Evangelos Kalogerakis, Stephen Giguere, and Thomas Funkhouser. Attribit:
698 content creation with semantic attributes. In *Proceedings of the 26th annual ACM symposium on*
699 *User interface software and technology*, pages 193–202, 2013.
- 700 [16] Adriana Schulz, Ariel Shamir, David IW Levin, Pitchaya Sitthi-Amorn, and Wojciech Matusik.
701 Design and fabrication by example. *ACM Transactions on Graphics (TOG)*, 33(4):1–11, 2014.
- 702 [17] Tianqiang Liu, Siddhartha Chaudhuri, Vladimir G Kim, Qixing Huang, Niloy J Mitra, and Thomas
703 Funkhouser. Creating consistent scene graphs using a probabilistic grammar. *ACM Transactions on*
704 *Graphics (TOG)*, 33(6):1–12, 2014.
- 705 [18] Minhyuk Sung, Hao Su, Vladimir G Kim, Siddhartha Chaudhuri, and Leonidas Guibas. Comple-
706 mentme: weakly-supervised component suggestions for 3D modeling. *ACM Transactions on Graphics*
707 *(TOG)*, 36(6):1–12, 2017.
- 708 [19] AJ Baddeley and MNM Van Lieshout. Stochastic geometry models in high-level vision. *Journal of*
709 *Applied Statistics*, 20(5-6):231–256, 1993.
- 710 [20] Xavier Descombes. Multiple objects detection in biological images using a marked point process
711 framework. *Methods*, 115:2–8, 2017.
- 712 [21] Xavier Descombes, Robert Minlos, and Elena Zhizhina. Object extraction using a stochastic birth-
713 and-death dynamics in continuum. *Journal of Mathematical Imaging and Vision*, 33(3):347–359,
714 2009.
- 715 [22] Mathias Ortner, Xavier Descombes, and Josiane Zerubia. Building outline extraction from digital
716 elevation models using marked point processes. *International Journal of Computer Vision*, 72(2):
717 107–132, 2007.
- 718 [23] Akos Utasi and Csaba Benedek. A 3-d marked point process model for multi-view people detection.
719 In *CVPR 2011*, pages 3385–3392. IEEE, 2011.
- 720 [24] Maria S Kulikova, Ian H Jermyn, Xavier Descombes, Elena Zhizhina, and Josiane Zerubia. Extraction
721 of arbitrarily-shaped objects using stochastic multiple birth-and-death dynamics and active contours.
722 In *Computational Imaging VIII*, volume 7533, page 753306. International Society for Optics and
723 Photonics, 2010.
- 724 [25] Dae Woo Kim, Camilo Aguilar, Huixi Zhao, and Mary L Comer. Narrow gap detection in microscope
725 images using marked point process modeling. *IEEE Transactions on Image Processing*, 28(10):5064–
726 5076, 2019.

- 727 [26] Quanhua Zhao, You Wu, Huabin Wang, and Yu Li. Road extraction from remote sensing image
728 based on marked point process with a structure mark library. *International Journal of Remote*
729 *Sensing*, 41(16):6183–6208, 2020.
- 730 [27] Amal Mbarki and Mohamed Naouai. A marked point process model for visual perceptual groups
731 extraction. In *2020 IEEE International Conference on Visual Communications and Image Processing*
732 *(VCIP)*, pages 511–514. IEEE, 2020.
- 733 [28] Xavier Descombes. *Stochastic geometry for image analysis*. John Wiley & Sons, 2013.
- 734 [29] John Quigley and Lesley Walls. A methodology for constructing subjective probability distributions
735 with data. In *Elicitation*, pages 141–170. Springer, 2018.
- 736 [30] John Quigley, Gavin Hardman, Tim Bedford, and Lesley Walls. Merging expert and empirical data
737 for rare event frequency estimation: Pool homogenisation for empirical bayes models. *Reliability*
738 *Engineering & System Safety*, 96(6):687–695, 2011.
- 739 [31] Gokula Vasantha, David Purves, John Quigley, Jonathan Corney, Andrew Sherlock, and Geevin
740 Randika. Common design structures and substitutable feature discovery in CAD databases. *Advanced*
741 *Engineering Informatics*, 48:101261, 2021. ISSN 1474-0346. doi: [https://doi.org/10.1016/j.aei.2021.](https://doi.org/10.1016/j.aei.2021.101261)
742 101261. URL <https://www.sciencedirect.com/science/article/pii/S1474034621000161>.
- 743 [32] John Quigley and Matthew Revie. Estimating the probability of rare events: addressing zero failure
744 data. *Risk Analysis: An International Journal*, 31(7):1120–1132, 2011.
- 745 [33] Harry Joe. Relative entropy measures of multivariate dependence. *Journal of the American Statistical*
746 *Association*, 84(405):157–164, 1989.
- 747 [34] Duncan Paterson and Johnathan Corney. Feature based search of 3d databases. In *International De-*
748 *sign Engineering Technical Conferences and Computers and Information in Engineering Conference*,
749 volume 50084, page V01BT02A010. American Society of Mechanical Engineers, 2016.
- 750 [35] Stuart Russell and Peter Norvig. *Artificial intelligence: a modern approach*. 2020.

$$\mathbb{E}[X] = \sum_{i=1}^{i_{max}} \sum_{m=1}^{m_{max}} \pi(i) p_{i,m} \left\{ \delta_{n_{im}} \frac{\sum_{j=1}^{n_{i,m}} x_{i,m,j}}{n_{i,m}} + (1 - \delta_{n_{im}}) \frac{1}{2} \right\} \quad (19)$$

$$\mathbb{E}[Y] = \sum_{i=1}^{i_{max}} \sum_{m=1}^{m_{max}} \pi(i) p_{i,m} \left\{ \delta_{n_{im}} \frac{\sum_{j=1}^{n_{i,m}} y_{i,m,j}}{n_{i,m}} + (1 - \delta_{n_{im}}) \frac{1}{2} \right\} \quad (20)$$

$$\mathbb{E}[XY] = \sum_{i=1}^{i_{max}} \sum_{m=1}^{m_{max}} \pi(i) p_{i,m} \left\{ \delta_{n_{im}} \frac{\sum_{j=1}^{n_{i,m}} x_{i,m,j} y_{i,m,j}}{n_{i,m}} + (1 - \delta_{n_{im}}) \frac{1}{4} \right\} \quad (21)$$

$$\mathbb{E}[X^2] = \sum_{i=1}^{i_{max}} \sum_{m=1}^{m_{max}} \pi(i) p_{i,m} \left\{ \delta_{n_{im}} \frac{\sum_{j=1}^{n_{i,m}} x_{i,m,j}^2}{n_{i,m}} + (1 - \delta_{n_{im}}) \frac{1}{4} \right\} \quad (22)$$

$$\mathbb{E}[Y^2] = \sum_{i=1}^{i_{max}} \sum_{m=1}^{m_{max}} \pi(i) p_{i,m} \left\{ \delta_{n_{im}} \frac{\sum_{j=1}^{n_{i,m}} y_{i,m,j}^2}{n_{i,m}} + (1 - \delta_{n_{im}}) \frac{1}{4} \right\} \quad (23)$$

753 S2 Mutual Information Distributions

754 The Mutual Information measure for the model is defined by ω in the following.

$$\omega = \mathbb{E} [\ln (f (X, Y, M))] - E [\ln (f (X))] - E [\ln (f (Y))] - E [\ln (f (M))]$$

755 Therefore we require the entropy measures for the joint and marginal distributions.

756 As we are investigating dependency within the historic data we use the following joint and marginal
757 probability density function where $\delta_{n_{i,m}}$ has been set to 1.

$$f(x, y, m) = \sum_{i=1}^{i_{\max}} \sum_{m=1}^{m_{\max}} \pi(i) \delta_m p_{i,m} \left(\frac{\sum_{j=1}^{n_{i,m}} \phi_j(x, y; \mu_{x,j} = x_{i,k,j}, \mu_{y,j} = y_{i,m,j}, \sigma)}{n_{i,m}} \right)$$

$$f(x) = \sum_{i=1}^{i_{\max}} \sum_{m=1}^{m_{\max}} \pi(i) p_{i,m} \frac{\sum_{j=1}^{n_{i,m}} g_j(x; \mu_{x,j} = x_{i,m,j}, \sigma)}{n_{i,m}}$$

$$f(y) = \sum_{i=1}^{i_{\max}} \sum_{m=1}^{m_{\max}} \pi(i) p_{i,m} \frac{\sum_{j=1}^{n_{i,m}} g_j(y; \mu_{y,j} = y_{i,m,j}, \sigma)}{n_{i,m}}$$

$$f(m) = \sum_{i=1}^{i_{\max}} \pi(i) p_{i,m}$$

758 Deriving the MI statistic will require either numerical integration of Monte Carlo simulations methods.

759 We propose the latter.

760 Assume we generate s' random simulations of (x, y) locations and feature m from $f(x, y, m)$. As we
761 increase the number of simulations we can obtain a more accurate estimate of the expectation.

$$\mathbb{E} [\ln (f (X, Y, M))] = \lim_{s' \rightarrow \infty} \frac{\sum_{s=1}^{s'} \ln (f (x_s, y_s, m_s))}{s'}$$

762 Similar arguments hold for the marginal on location (x, y) . For feature direct calculation of the
763 entropy is straightforward.

$$\mathbb{E} [\ln (f (M))] = \sum_{i=1}^{i_{\max}} \sum_{m=1}^{m_{\max}} \pi(i) p_{i,m} \ln \left(\sum_{i=1}^{i_{\max}} \pi(i) p_{i,m} \right)$$

764 **S3 Case Study Demonstration**

765 **S3.1 Feature Location Probability Function Assessment**

766 For illustration, we step through an example using the data on the $x = 0$ face, and focus on one test
767 design, and consider it a new design. With no features added to this on the $x = 0$ face, a uniform
768 prior was assumed on which design from the training set the new design is most similar to, and use this
769 to update the predictive distribution. Figure 10 provides a visualization of the predictive density. The
770 darker red regions indicate areas of higher feature probability which reflects the contents on the training
771 dataset with the central bore hole most prevalent. The blue numbers are used to indicate the positions
772 of the hole features in the test design.

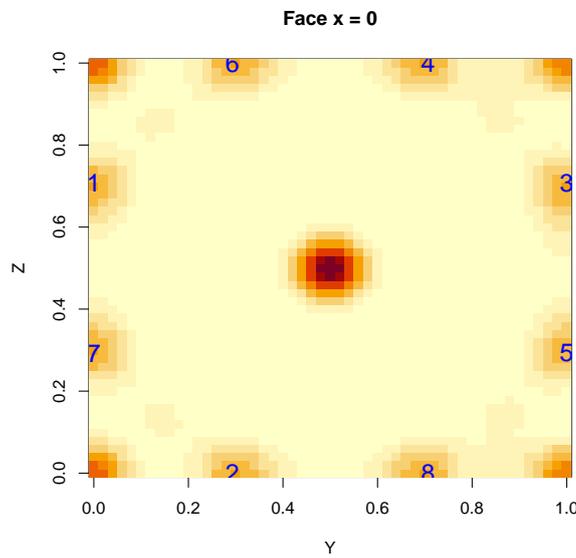


Figure 10: The predictive density for a new design given that no features have been added to the face with the blue numbers indicating the positions of the features that will be added showing strong agreement between the initial prediction and realization.

773 One feature is then added at the position observed in the test design. Given this additional informa-
774 tion, the prior on which design from the training dataset the new design is most similar to is updated.
775 Figure 11 illustrates how the prior on each training set design changes given this one observation; the
776 red dashed line indicates the Uniform prior of no information and the probability of test designs either
777 increase or decrease given the addition of the feature to the new design. Importantly, while the probab-
778 ilities can be very small they are all non-zero. Moreover, we see an clear indication of a subset of designs
779 types that are more likely candidates for the design being constructed.

780 **S3.2 Predictive Feature Location Function Predictive Performance**

781 This posterior is then used to update the predictive distribution, with the density, aggregated over all
782 features, shown in Fig. 12; the green triangle indicates the position of the feature that was added to the

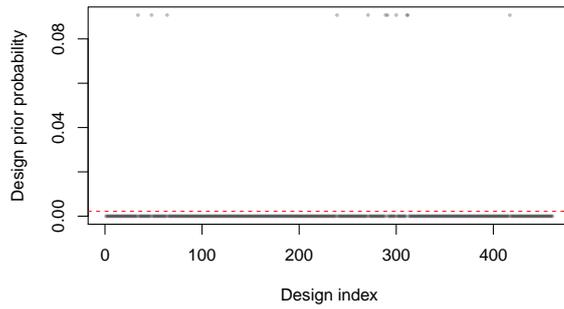


Figure 11: The posterior probability on which design from the training dataset a new design is most similar to. The red dashed line indicates the prior Uniform probability where all designs have equal probability.

783 new design. The addition of one feature shifts the probability to similar designs in the training dataset,
 784 and which captures the similar scaled coordinates.

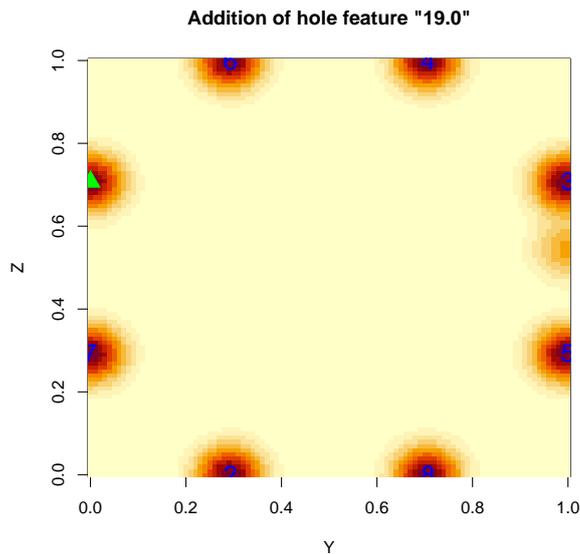


Figure 12: The updated predictive density for a new design given one feature, hole diameter “18.0”, has been added at the position marked by the green triangle, showing stronger agreement between the prediction and the realizations than with the prior distribution.

785 Additionally, due to how the density is generated, the predictions can be presented at the feature
 786 level, and the ranked probability of feature inclusion calculated at each spatial coordinate. See Fig. 13
 787 for the predictive distribution for the four features with the highest probability aggregated over the unit
 788 square; 94% of the predictive density weight is on the “19.0” feature, 5% on “6.4”, and all other features
 789 having lower than 0.1% (note that the density color scale is not transferable between subplots i.e. the
 790 dark red may represent lower density between images).

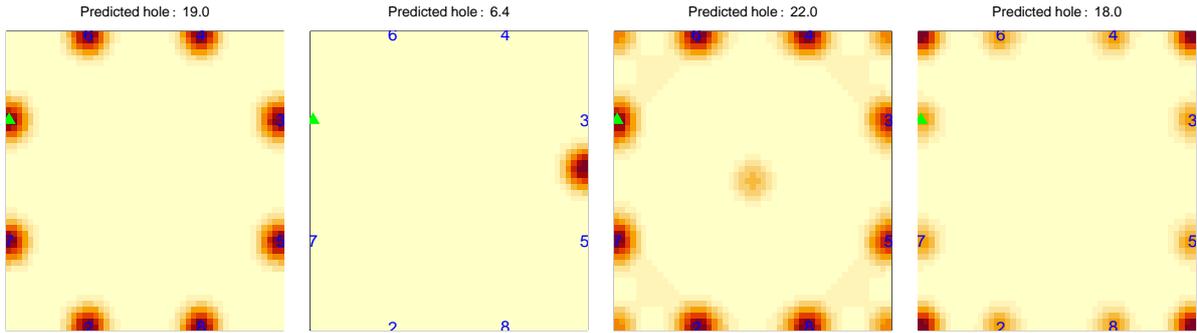


Figure 13: The updated predictive density for four different features for a new design given one feature, hole diameter “18.0”, has been added at the position marked by the green triangle

791 S3.3 Measuring Predictive Performance

792 Fig. 12 also illustrates how the system’s predictive performance was evaluated. Using the test design,
 793 the observed features can be sequentially added to a new design at the observed test positions one at
 794 a time. A distance measure was then calculated after each addition between the positions of each of
 795 the remaining holes in the test design and the nearest predictive mode – the distance from the observed
 796 features in blue to the areas of high density in red – e.g. how well does the predictive density capture
 797 the feature positions. Details of how the distances were calculated are provided in S4.

798 A score was also generated to measure how accurately a hole feature is predicted at a specific position.
 799 The predicted density of each feature at the position of an observed hole in the test design were ranked,
 800 and this ranked list was then compared to the feature observed in the test design at that position;
 801 reciprocal rank was used as the performance measure.

802 These measures were calculated for all designs in the test datasets. The distances to remaining hole
 803 features and the ranks were averaged within each fold given how many holes have been added to the new
 804 design.

805 **S4 Evaluating Distance from Feature to Nearest Predictive Mode**

806 A hill-climbing algorithm [35] was used to calculate the distance from the remaining features in the
807 design to the nearest predictive mode; starting from the feature position moves in steps of one around
808 the predicted intensity grid were evaluated until there were no increasing moves. The Euclidean distance
809 between the feature position and the mode was then calculated. Euclidean distance was measured in two
810 scales; one which accounts for the difference in scaling the dimensions to $[0, 1]$ and another in the raw
811 counts of moving from one grid position to the next. For example in the Fig. 14, the predictive density
812 on the $x = 0$ is calculated at discrete points in y and z , shown by the black points. We observe a point
813 from the test design, in red. The search for the nearest predictive mode evaluates all of the adjacent grid
814 positions to the current position in the search by comparing the density (dependent on scaled distance)
815 of the neighboring discrete points to the observed; color blue. If a neighboring grid point has a higher
816 value (green) then the search moves to that position. This continues until there are no increasing moves
817 – this position is then taken as the nearest local mode to the observed point.

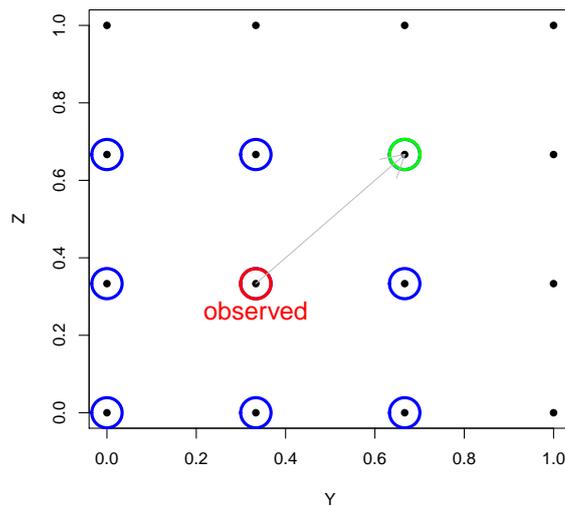


Figure 14: Illustration of hill-climbing search to nearest predictive mode. The $[0, 1]$ interval is divided into a N by N grid, in this case $N = 4$. The kernel density is estimated at the discrete black points. A feature is observed at a specific position, shown in red, and the mode is located by moving in increasing steps between neighboring grid positions.

818 **S5 Rotation of Features in Two Dimensions**

819 The following example illustrates how the features of a part can be rotated to be more in alignment with
820 the features in the database of designs. For an example, a part may have been designed around a rotated
821 axis relative to the other parts in the database. Figure 15 shows the part from earlier rotated 20 degrees
822 on the $x = 0$ face. The black points show the superimposition of all features instances from the training
823 set on the $x = 0$ face – the average design, the blue triangles illustrate the positions of the features from a
824 misaligned design – the design to be rotated – and the green stars give the feature positions post rotation.
825 The rotation was seen to improve predictive ability; the average unscaled distance (grid points) from the
826 features to the predicted modes was seven steps before rotation but one after.

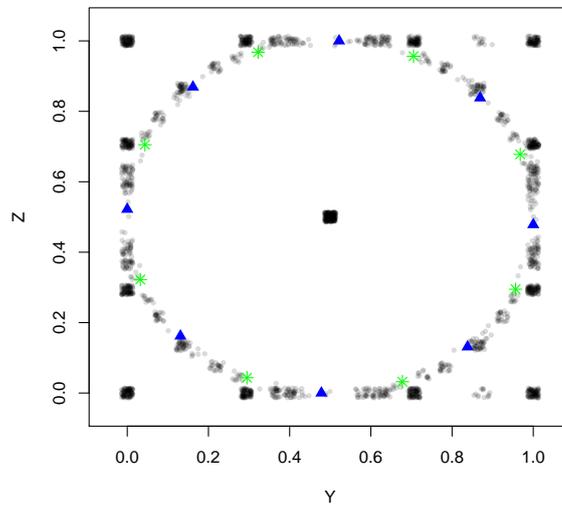


Figure 15: A superimposition of all features instances from the training set on the $x = 0$ face. The blue triangles illustrate the positions of the features from a misaligned design. The green stars give their positions post rotation

828 **S6.1 Different Spatial Orientation of Features**

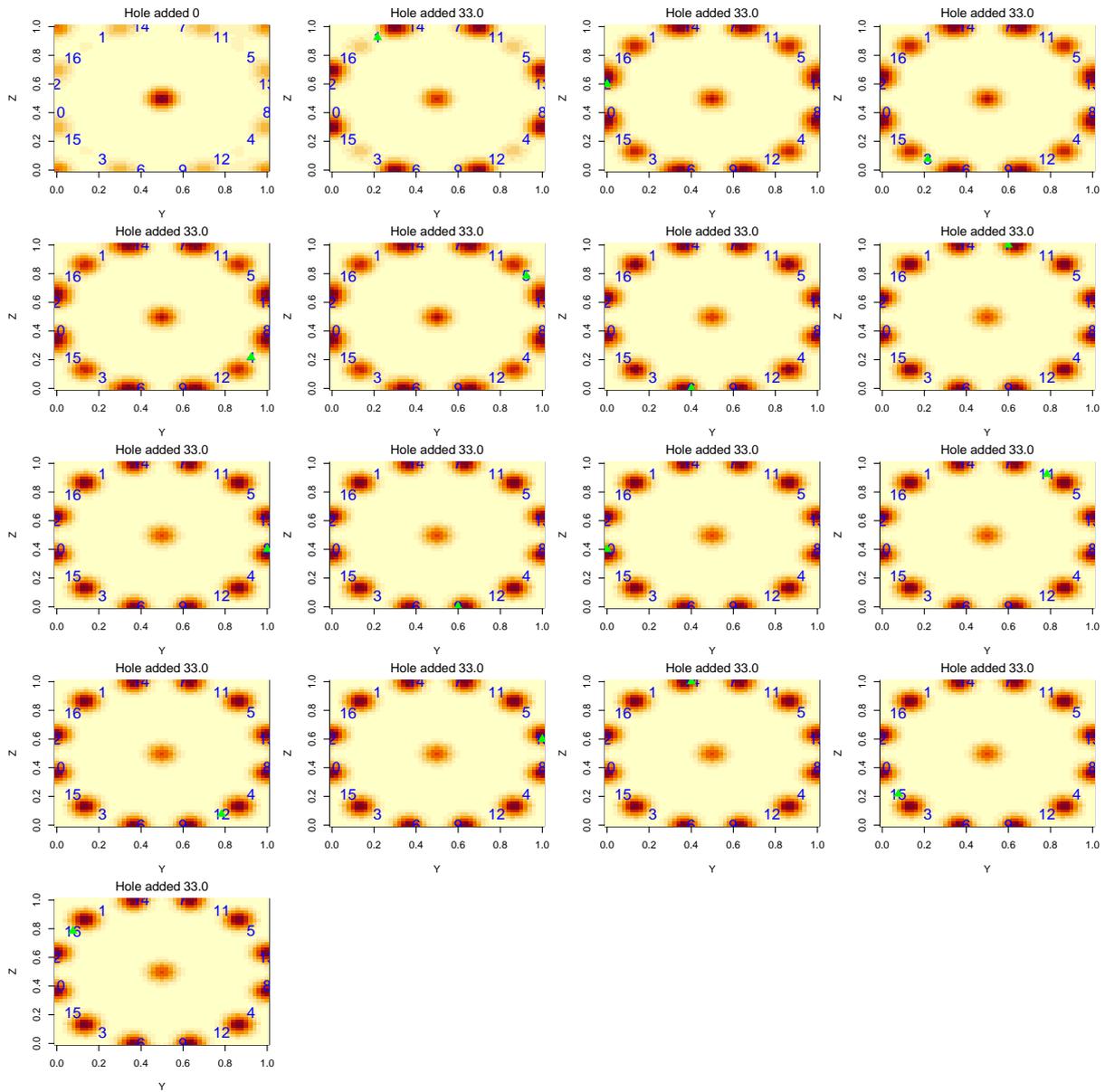


Figure 16: Predictive density for a new (test) design with greater feature occurrence than is present in the training data designs

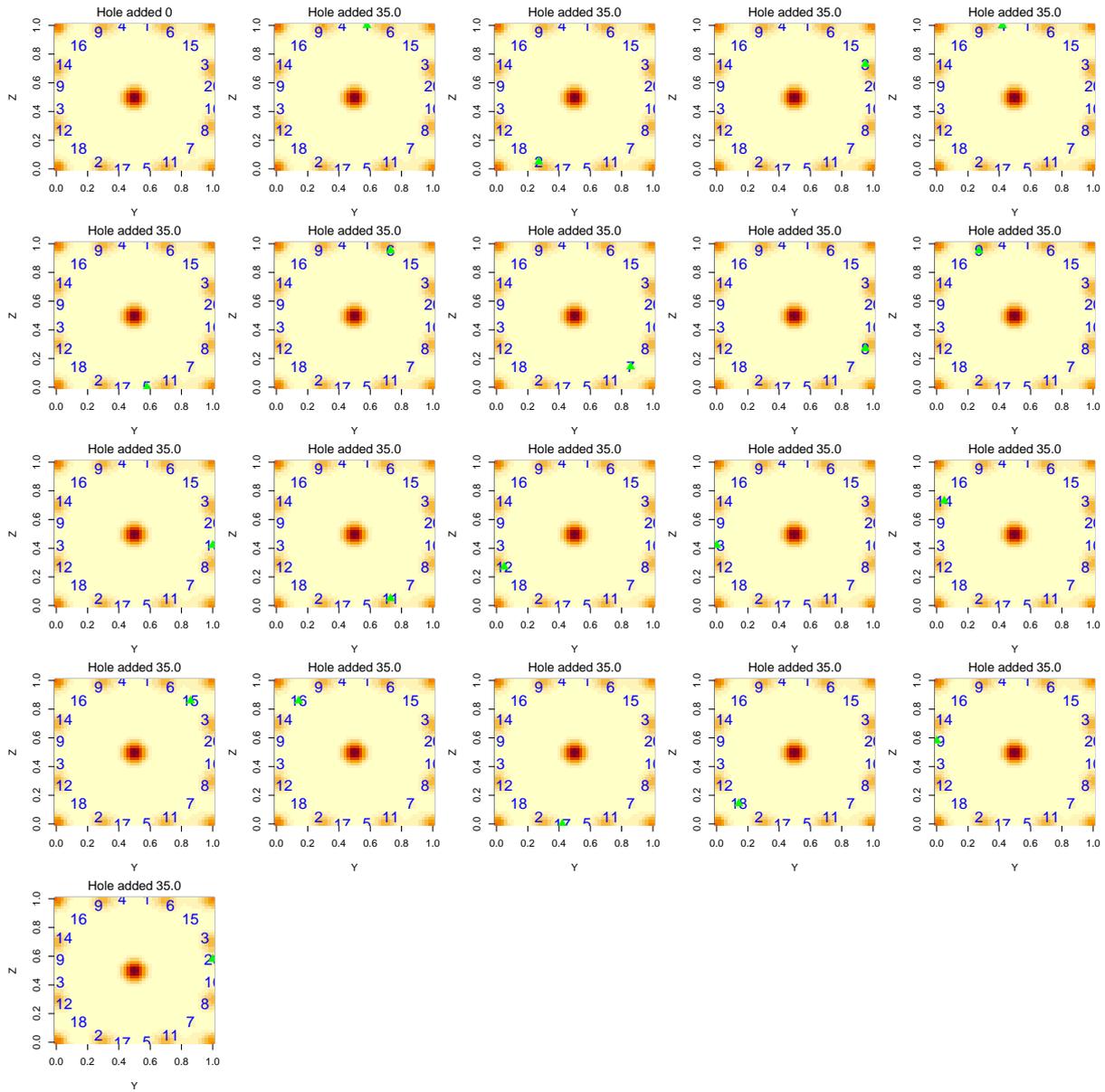


Figure 17: Predictive density for a new (test) design in which a hole feature is used for a different purpose than in the training data designs

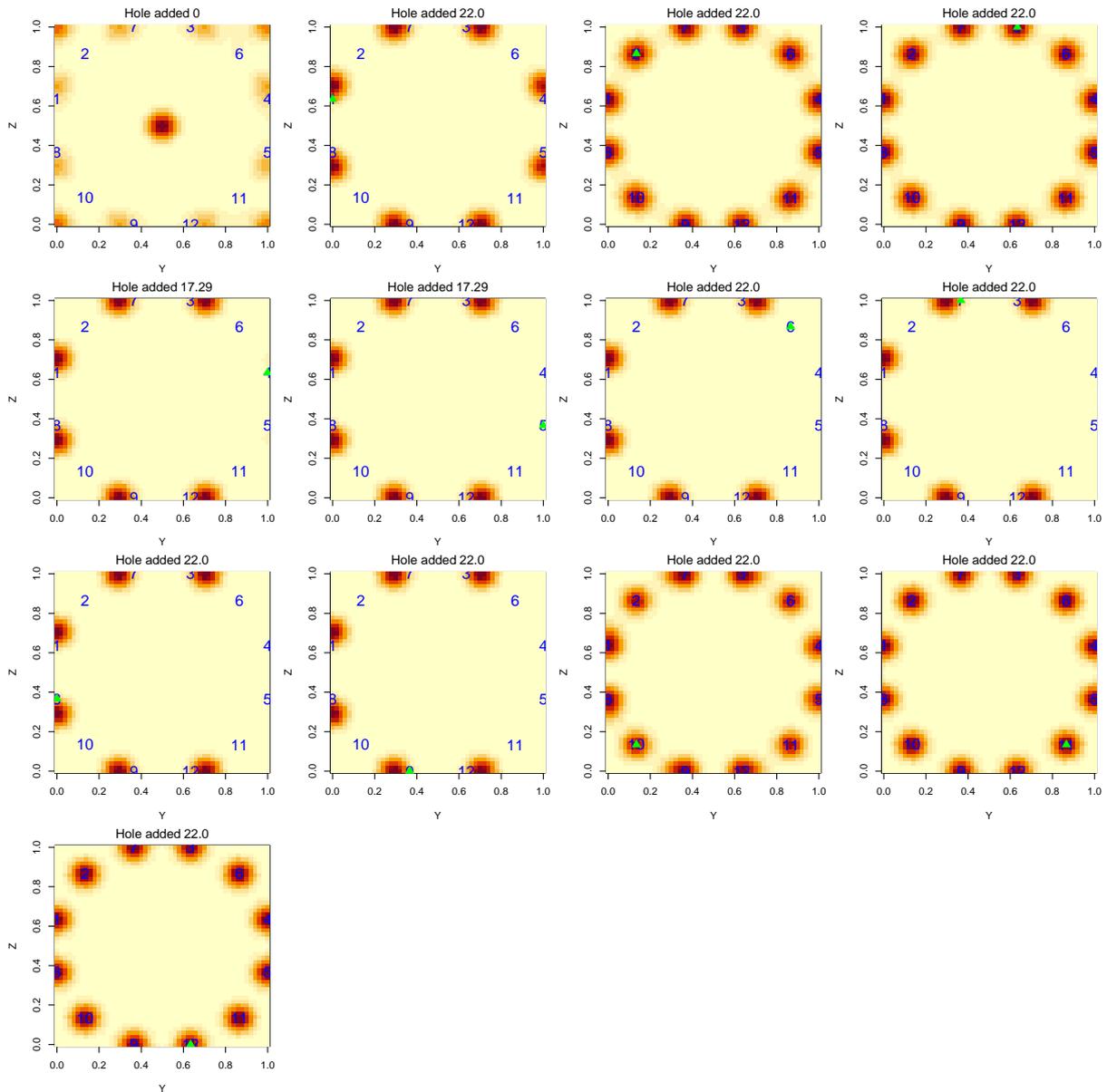


Figure 18: The order that hole features are added to a new design impacts the predictive density and thus decision support.

831 **S7 Additional Modeling Results**

832 **S7.1 Prediction Results for the $y = 0$ Face**

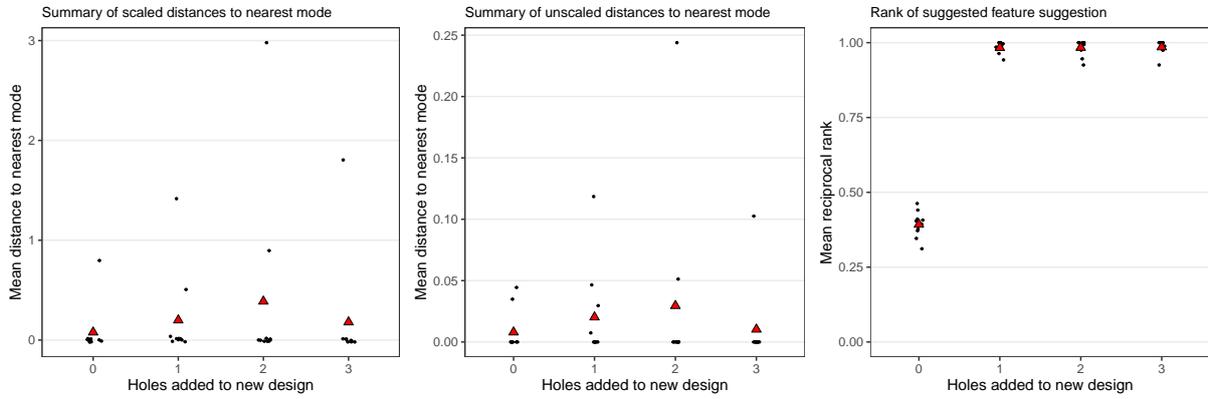


Figure 19: Cross-validation predictive performance. The x-axis indicates how many holes have been added to a new design and the y-axis records the average distance or rank. The black points give the within fold average measure and the red triangle the average performance across the folds.

833 **S7.2 Prediction Results for the $z = 0$ Face**

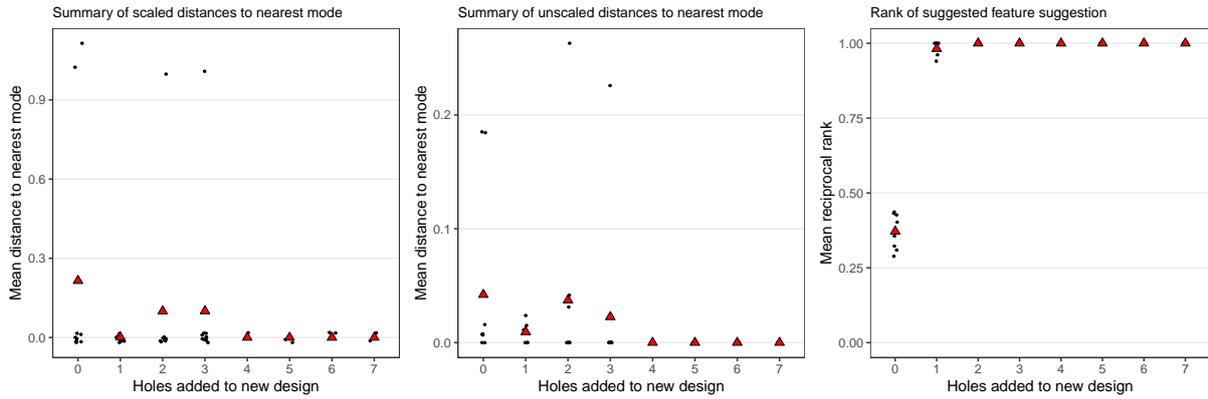


Figure 20: Cross-validation predictive performance. The x-axis indicates how many holes have been added to a new design and the y-axis records the average distance or rank. The black points give the within fold average measure and the red triangle the average performance across the folds.

834 List of Figures

835	1	User Interface for Assembly-based modeling using relative attributes [15]	7
836	2	Input scene mapped with labels and converted into the hierarchical form using probabilistic	
837		grammar [17]	8
838	3	System schematic for feature prediction through Bayesian updating.	11
839	4	Example of a valve body from CAD design database. Figure 4a shows an image of a valve	
840		design and Fig. 4b the scaled positions of the extracted hole features. The different colors	
841		are used to identify the different diameters of the holes	20
842	5	Superimposition of scaled hole feature coordinates for all the data on the specified face of	
843		the cube.	20
844	6	Cross-validation predictive performance. The x-axis indicates how many holes have been	
845		added to a new design and the y-axis records the average distance or rank. The black	
846		points give the within fold average measure and the red triangle the average performance	
847		across the folds.	22
848	7	Cross-validation predictive performance after omitting unusual designs from the test set.	
849		The x-axis indicates how many holes have been added to a new design and the y-axis records	
850		the average distance or rank. The black points give the within fold average measure and	
851		the red triangle the average performance across the folds.	23
852	8	Predictive density calculated under a Normal kernel using different standard deviations.	
853		The colourbar legend describes how the image colour maps to the data, with dark red	
854		indicating regions of higher predictive density	25
855	9	The feature predictions were mapped onto two incomplete designs. Regions of higher	
856		predictive density are illustrated by dark red. From left to right, the plots show the	
857		complete part, the predictive distribution when no features have been added, and so is	
858		completely informed by the training data, and lastly the predictive distribution after the	
859		addition of a central bore hole and one bolt hole feature. The updated distribution displays	
860		the increased level of belief in the positions of subsequent features, with areas of higher	
861		density and lower variance. The black circle in (a) is used to indicate the prediction	
862		region which has been truncated and the red and blue circles in (b) are used to highlight	
863		the changes in the predictive density once a choice has been made. The green sections	
864		represent areas outside the normalized prediction region.	27
865	10	The predictive density for a new design given that no features have been added to the face	
866		with the blue numbers indicating the positions of the features that will be added showing	
867		strong agreement between the initial prediction and realization.	35

868	11	The posterior probability on which design from the training dataset a new design is most similar to. The red dashed line indicates the prior Uniform probability where all designs have equal probability.	36
869			
870			
871	12	The updated predictive density for a new design given one feature, hole diameter “18.0”, has been added at the position marked by the green triangle, showing stronger agreement between the prediction and the realizations than with the prior distribution.	36
872			
873			
874	13	The updated predictive density for four different features for a new design given one feature, hole diameter “18.0”, has been added at the position marked by the green triangle	37
875			
876	14	Illustration of hill-climbing search to nearest predictive mode. The $[0, 1]$ interval is divided into a N by N grid, in this case $N = 4$. The kernel density is estimated at the discrete black points. A feature is observed at a specific position, shown in red, and the mode is located by moving in increasing steps between neighboring grid positions.	38
877			
878			
879			
880	15	A superimposition of all features instances from the training set on the $x = 0$ face. The blue triangles illustrate the positions of the features from a misaligned design. The green stars give their positions post rotation	39
881			
882			
883	16	Predictive density for a new (test) design with greater feature occurrence than is present in the training data designs	40
884			
885	17	Predictive density for a new (test) design in which a hole feature is used for a different purpose than in the training data designs	41
886			
887	18	The order that hole features are added to a new design impacts the predictive density and thus decision support.	42
888			
889	19	Cross-validation predictive performance. The x-axis indicates how many holes have been added to a new design and the y-axis records the average distance or rank. The black points give the within fold average measure and the red triangle the average performance across the folds.	43
890			
891			
892			
893	20	Cross-validation predictive performance. The x-axis indicates how many holes have been added to a new design and the y-axis records the average distance or rank. The black points give the within fold average measure and the red triangle the average performance across the folds.	43
894			
895			
896			