

基于变分自编码器的实验设计^①

张志博^{1,2}, 康达周^{1,2,3}

¹(南京航空航天大学 计算机科学与技术学院/人工智能学院, 南京 211106)

²(南京航空航天大学 高安全系统的软件开发与验证技术工信部重点实验室, 南京 211106)

³(软件新技术与产业化协同创新中心, 南京 210023)

通信作者: 康达周, E-mail: dzkang@nuaa.edu.cn



摘要: 针对现有实验设计方法难以对复杂系统进行高效实验设计的问题, 本文提出了一种基于变分自编码器的实验设计方法, 首先利用实验历史记录数据训练变分自编码器将复杂的实验样本空间编码到一个较为简单的隐变量空间, 然后在该隐变量空间里进行取样, 最后通过解码器还原产生新的实验样本, 完成实验设计. 通过对比本文方法与数种基准实验设计方法的结果在拟合直航鱼雷命中模型时的表现情况, 表明在取相同样本数的情况下, 本文方法可以优化实验设计, 提高实验效率.

关键词: 复杂系统; 实验设计; 变分自编码器; 支持向量回归

引用格式: 张志博, 康达周. 基于变分自编码器的实验设计. 计算机系统应用, 2022, 31(3): 113-121. <http://www.c-s-a.org.cn/1003-3254/8333.html>

Design of Experiments Based on Variational Auto-encoder

ZHANG Zhi-Bo^{1,2}, KANG Da-Zhou^{1,2,3}

¹(College of Computer Science and Technology/College of Artificial Intelligence, Nanjing University of Aeronautics and Astronautics, Nanjing 211106, China)

²(Key Laboratory of Safety-Critical Software, Ministry of Industry and Information Technology, Nanjing University of Aeronautics and Astronautics, Nanjing 211106, China)

³(Collaborative Innovation Center of Novel Software Technology and Industrialization, Nanjing 210023, China)

Abstract: Given that the existing experiment methods are unable to perform the efficient design of experiments for complex systems, this study proposes a design of experiments method based on the variational auto-encoder. First, experimental historical data are used to train the variational auto-encoder to encode the complex experimental sample space into a relatively simple latent variable space. Then, samples are obtained from the latent variable space. Finally, new experimental samples are generated by the decoder through restoration, and the design of experiments is achieved. The performance of the proposed method in fitting the hit model of the straight-running torpedo is compared with those of several benchmark design of experiments methods. It is shown that with the same number of samples, the proposed method can optimize the design of experiments and improve the efficiency of the experiments.

Key words: complex system; design of experiments; variational auto-encoder (VAE); support vector regression (SVR)

复杂系统^[1]是一种包含很多相互作用子系统的系统. 它具有非线性、高不确定性、高维度、多层次以及高相关性的特点. 空气动力系统^[2]、股票市场系统^[3]、暴雨洪涝人口风险系统^[4]就是一些的复杂系统例子.

为研究自然、工程和社会科学中的现象与原理, 探究这些复杂系统的内在性质, 人们对复杂系统建模并利用计算机对它们进行仿真与实验, 并让实验能够尽可能地接近真实世界, 不断向着复杂化发展. 比如, 在武

① 基金项目: 十三五装备预研共用技术项目 (41402020101); 基础科研项目 (JCKY2020605C003)

收稿时间: 2021-05-05; 修改时间: 2021-05-19; 采用时间: 2021-06-07; csa 在线出版时间: 2022-01-24

器装备领域,为了测定直航鱼雷^[5]在不同复杂想定下的打击范围和命中能力,需要做大量的仿真实验.为了让实验能在较小的代价下执行,需要进行高效实验设计.但复杂系统的特点会给复杂系统实验设计带来非线性、相关性、不确定性和规模性等问题^[6],传统的实验设计方法只能部分解决这些问题,例如:完全析因设计^[7]适用于非线性系统,但需要的实验样本和实验次数非常多;响应面设计^[8]用于分析系统中因子间的相互作用,但当因子数量增多时,计算量会急剧增加;为了降低实验中的不确定性,可以通过大量的随机重复抽样设计,但这同样的会极大地增加实验规模.以上问题都会使得复杂系统实验设计最终走向规模性,即产生维度灾难^[9].正交设计^[10]、均匀设计^[11]、拉丁超立方设计^[12]这类遵循“充满空间”性质的实验设计方法虽然可以通过优化实验取样来缩减实验规模,但仅在均匀空间中最有效,对复杂系统实验设计效率的提升有限.因此,迫切需要一种适用于复杂系统、可以大幅提高实验效率的实验设计方法.

变分自编码器 (variational auto-encoder, VAE)^[13]是一种借鉴变分推断原理的深度神经网络,凭借其优秀的性能,近年来在图像处理^[14]、语音识别^[15]、文本生成^[16]等应用领域取得了巨大的成功.变分自编码器作为一种生成模型,含有多个隐藏层,包括了编码器、解码器和隐变量空间.编码器可以将复杂分布的样本数据投影到简单的隐变量空间,实现对样本数据的特征提取.而解码器可以通过还原在隐变量空间里取样的数据,生成新的样本.和一般生成模型只考虑重构损失函数不同的是,变分自编码器同时考虑了重构损失函数和 KL 损失函数,使得生成的样本具有更强的泛化能力,可以提高样本空间取点的覆盖率,这给复杂系统实验设计带来了新的思路:一方面,如图 1 所示,将复杂系统实验的输入样本和输出响应结果一同作为实验历史记录数据训练变分自编码器,通过生成新的实验样本,综合考虑复杂系统实验输入样本 x_i 的分布、输出结果 y_i 的分布和它们之间的对应关系 $y_i = f(x_i)$ 的分布,使得 3 个分布都尽可能地“充满空间”,实现更加高效的实验设计.另一方面,VAE 将复杂分布的输入实验样本投影到的隐变量空间会尽可能满足标准正态分布,而在符合标准正态分布的空间里进行实验设计有 3 个好处.

- (1) 通过适应随机采样,防止出现过拟合;
- (2) 隐变量空间具有良好的分布性质,易于抽样

处理;

(3) 生成的样本点是在隐变量空间采样得到,可以通过调整采样来进行优化.

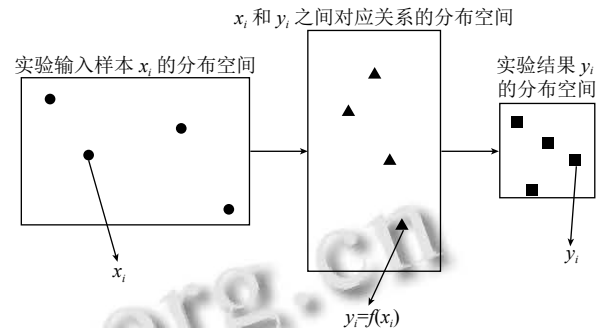


图 1 提高实验设计效率的思路是使这 3 个分布尽可能充满空间

基于上述分析,本文提出一种基于 VAE 的实验设计方法,将复杂系统实验历史记录数据投影到一个满足正态分布的隐变量空间,并通过对隐变量空间采样生成新的实验设计方案.实验结果表明,本文实验设计方法系统性地解决了复杂系统自身特性导致的实验设计规模性问题.

本文接下来的内容安排如下:第 1 节详细阐述了本文的方法;第 2 节介绍了实验有效性的评价方法;第 3 节通过一个直航鱼雷命中模型的案例对本文方法进行了验证并展示与分析了实验结果;第 4 节总结了全文并得出了结论.

1 一种基于变分自编码器的复杂系统实验设计方法

本节首先形式化定义了待解决问题,然后详细介绍了本文方法的总体架构和各模块内容.

1.1 问题定义

复杂系统的特性可以表示为其输入到输出的一个函数: $\mathbf{y} = f(\mathbf{x}) = f(x_1, x_2, \dots, x_d)$, 其中, $\mathbf{x} = (x_1, x_2, \dots, x_d)$ 是输入变量的向量, d 是因子的数目, $\mathbf{y} = (y_1, y_2, \dots, y_{d_y})$ 为该复杂系统对输入 \mathbf{x} 的响应, d_y 是响应的数目.

实验设计目的之一是提高实验设计取样的有效性,即在相同实验样本数量的情况下,提高选取样本的准确性.对于待比较的不同实验设计方法,令每种实验设计方法产生 k 个实验样本,可以通过这 k 个样本所拟合出的模拟系统 $\mathbf{y}' = f'(\mathbf{x})$ 与真实系统 $f(\mathbf{x})$ 的吻合程度来评价它们的实验准确性. $f'(\mathbf{x})$ 与 $f(\mathbf{x})$ 的吻合程度可以

用均方根误差 ($RMS E$)、最大残差 (MRE)、平均绝对误差 (MAE)、 R^2 决定系数表征:

$$RMS E = \sqrt{\frac{\sum_{i=1}^k (y_i - y'_i)^2}{k}} \quad (1)$$

$$MRE = \max(|y_i - y'_i|) \quad (2)$$

$$MAE = \frac{\sum_{i=1}^k |y_i - y'_i|}{k} \quad (3)$$

$$R^2 = 1 - \frac{\sum_{i=1}^k (y_i - y'_i)^2}{\sum_{i=1}^k (y_i - \bar{y})^2} \quad (4)$$

在式 (4) 中, \bar{y} 是复杂系统 k 个真实响应 y_i 的平均值. $RMS E$ 、 MRE 、 MAE 越小, R^2 决定系数越接近 1, 意味着 $f'(x)$ 与 $f(x)$ 的吻合程度越高、越相关, 即实验取点越有效. 因此, 待解决的问题可以表示成: 找到一种实验设计方法 E , 用该方法生成的 k 个实验样本拟合复杂系统 $f(x)$ 的模拟系统 $f'(x)$, 并可以使得 $f'(x)$ 与 $f(x)$ 之间尽可能地吻合.

1.2 总体架构

基于变分自编码器的实验设计方法的总体架构如图 2 所示, 它包含了 3 个模块.

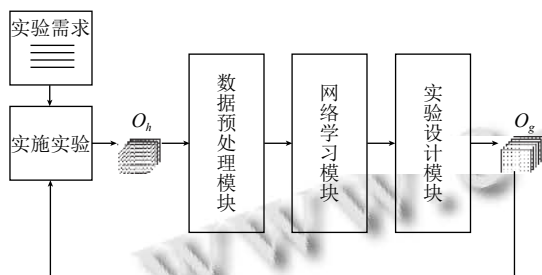


图 2 总体架构

实验需求确定了实验的背景、对象、参数以及属性的取值范围和概率分布, 为本实验设计方法的输入. 实验设计方案即实验样本集 O_g 为本实验设计方法的输出. 实验历史记录集 O_h 用于训练本实验设计方法的变分自编码器.

基于变分自编码器的实验设计方法各个模块的作用和总体步骤如下.

(1) 依据实验需求, 设计并实施实验, 得到初次的

实验历史记录集.

(2) 数据预处理模块将实验历史记录集转化成变分自编码器可以使用的实验历史记录向量集.

(3) 网络学习模块将此实验历史记录向量集用作训练集, 训练变分自编码器.

(4) 实验设计模块在训练好的变分自编码器的隐变量空间里采样, 利用解码器还原采样点生成实验样本, 输出复杂系统实验设计方案.

(5) 对输出的复杂系统实验设计方案实施实验, 并将得到的实验历史记录加入原先的实验历史记录集, 得到新的实验历史记录集, 然后回到步骤 (2) 进行下一轮的实验设计, 使得实验设计取样得到进一步的优化.

1.3 数据预处理模块

变分自编码器的训练集需要向量化或数值化的样本, 但是实验历史记录并不一定都满足这个条件. 在训练变分自编码器的过程中, 为了能让损失函数更容易收敛到最佳结果, 需要防止样本中有因子取值范围过大或存在极端值的情况. 为此需要对实验历史记录进行数据预处理.

对于实验历史记录中无序离散值的属性来说, 需要对其进行独热编码. 对于有序离散值的属性需要对其进行标签编码.

而对于实验历史记录中连续数值, 防止取值范围过大或有极端值情况的策略是使用 Minmax 归一化:

$$\text{Minmax}(O) = \frac{X_i - X_{\min}}{X_{\max} - X_{\min}} \quad (5)$$

其中, O 是实验历史记录集, 实验历史记录记录 $X_i \in O$, X_{\max} 是实验历史记录最大值, X_{\min} 是实验历史记录最小值.

因此, 在数据预处理模块, 对于复杂系统实验历史记录集 $O_h = ((x_1, y_1), (x_2, y_2), \dots, (x_n, y_n))$ 的每一个样本 (x_i, y_i) 对它进行如下步骤:

(1) 判断实验历史记录每个因子的类型.

1) 对实验历史记录中的有序离散取值的属性进行标签编码;

2) 对实验历史记录中的无序离散取值的属性进行独热编码;

3) 对实验历史记录中的连续数值的属性进行归一化.

(2) 对步骤 (1) 中得到的结果进行向量组装, 得到数据预处理后的实验历史记录.

数据预处理后的实验历史记录记为 $e_i = (e_1, e_2, \dots, e_d) \in O_e$, O_e 为数据预处理后的实验历史记录集. 通常

数据预处理后样本的维度不会低于数据预处理前实验历史记录的维度, 即 $d' \geq d + d_y$.

本文方法用 Python 库 SciKit-Learn^[17] 的 preprocessing 模块的 OneHotEncoder 方法、LabelEncoder 方法和 MinMaxScaler 方法实现数据的独热编码、标签编码和归一化。

1.4 网络学习模块

网络学习模块的主体部分是变分自编码器。

变分自编码器由编码器 $q_\phi(z|\mathbf{x})$ 、解码器 $p_\theta(z|\mathbf{x})p_\theta(\hat{\mathbf{x}}|z)$ 和一组隐变量 z 组成, \mathbf{x} 是输入的样本, $\hat{\mathbf{x}}$ 是重构或生成的样本。而编码器由均值编码器和方差编码器组成。均值编码器产生每个样本专属的在隐变量空间中的均值 μ_i , 方差编码器产生每个样本专属的在隐变量空间中的方差 σ_i^2 。从 $\varepsilon_i \sim N(0,1)$ 中采样, 并通过 $z_i = \mu_i + \sigma_i \cdot \varepsilon_i$ 的重参数技巧, 即可完成对隐变量 z 的采样, 经过此法采样的 z , 满足 $p_\theta(z) \sim N(0,1)$ 。

编码器和真实后验概率分布 $p_\theta(z|\mathbf{x})$ 之间的相似度

用 KL 散度衡量:

$$KL(q_\phi(z|\mathbf{x}) \parallel p_\theta(z|\mathbf{x})) = E_{q_\phi(z|\mathbf{x})} \log \frac{q_\phi(z|\mathbf{x})}{p_\theta(z|\mathbf{x})} \quad (6)$$

整个变分自编码器训练就是要去优化编码器模型的参数 ϕ 和解码器模型的参数 θ , 使得学习的后验概率分布和真实后验概率分布尽可能地相似, 即最小化它们之间的 KL 散度:

$$\phi, \theta = \arg \min_{\phi, \theta} KL(q_\phi(z|\mathbf{x}) \parallel p_\theta(z|\mathbf{x})) \quad (7)$$

定义生成模型 $p_\theta(\hat{\mathbf{x}}|z)$ 服从伯努利分布, 由文献 [18] 的数学推导, 得到变分自编码器的损失函数为:

$$L(\theta, \phi) = \underbrace{KL(q_\phi(z|\mathbf{x}) \parallel p_\theta(z))}_{KL\text{-loss}} - \underbrace{E_{q_\phi(z|\mathbf{x})} \log p_\theta(\hat{\mathbf{x}}|z)}_{\text{Reconstruction-loss}} \quad (8)$$

其中, 第一项是变分自编码器的 KL 损失函数, 第二项是变分自编码器的重构损失函数。

当完成对变分自编码器的训练后, 就可使用解码器生成新的样本。网络学习模块的结构如图 3 所示。

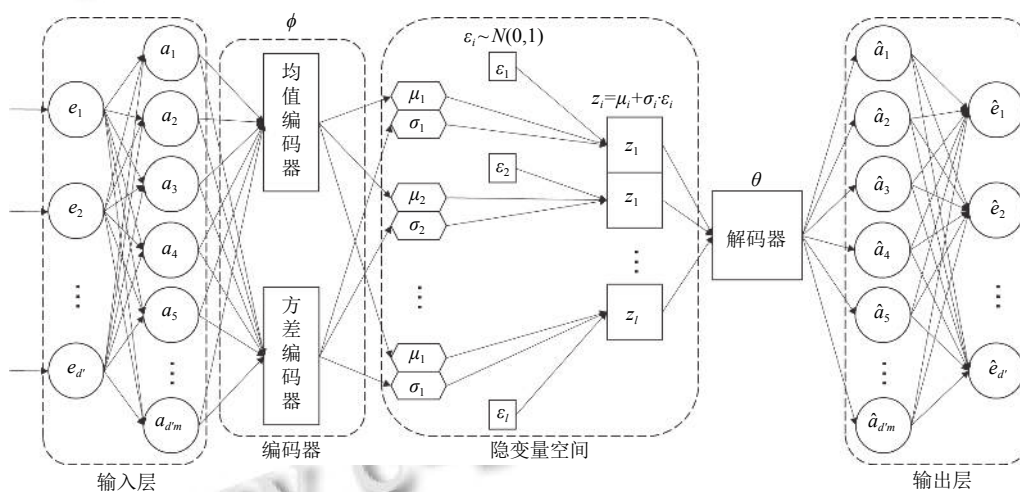


图 3 网络学习模块结构

由于当输入样本的维度不高时, 使用深度生成模型变分自编码器不一定能很好地学习出样本的内在性质, 为了获得更精确的学习效果, 需要酌情增加输入样本的维度。网络学习模块在常规变分自编码器结构中增加了输入层和输出层。输入层的作用是将输入样本的维度扩大 m 倍, 与输入层相对应的输出层的作用是将维度扩大 m 倍的样本再还原成原本维度的样本。

网络学习模块的输入是经过数据预处理的实验历史记录集 O_e 中的样本 e_i , $\mathbf{a}_i = (a_1, a_2, \dots, a_{d'm})$ 是 e_i 经过输入层的样本。 l 为隐变量空间的维度, 通常情况下

$d'm \gg l$ 。网络学习模块的输出是编码器的参数和解码器的参数。

1.5 实验设计模块

实验设计模块在由网络学习模块求出的隐变量空间里做抽样实验设计, 将原本分布不均的复杂系统实验历史记录空间投影到低维的服从标准正态分布的隐变量空间。因为训练中引入了随机的扰动, 这个过程可以确保与原始实验样本编码的潜在位置靠近的每个点都能被解码为与原始实验历史记录类似的实验历史记录, 从而迫使隐变量空间能够连续地有意义。隐变量空

间中任意两个相邻的点都会被解码为高度相似的实验历史记录. 隐变量空间的连续性以及低维度, 使得隐变

量空间非常适合进行抽样. 实验设计模块结构如图4所示.

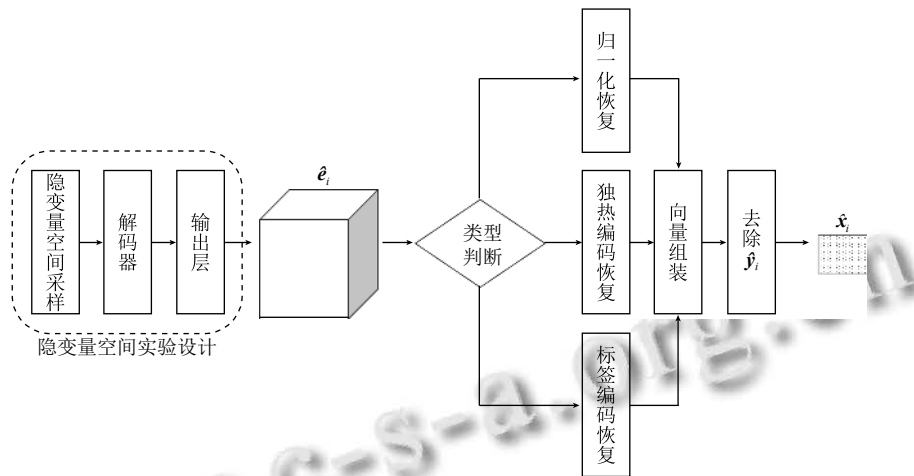


图4 实验设计模块结构

实验设计模块在隐变量空间里做实验设计, 相当于是对隐变量 z 采样, 让 z 按标准正态分布随机抽样. 完成隐变量空间的采样后, 隐变量 z 通过解码器生成了具有 $d'm$ 维的向量 $\hat{a}_i = (\hat{a}_{11}, \hat{a}_{12}, \dots, \hat{a}_{d'm})$. 再经过输出层的样本维度还原生成了和网络学习模块输入向量有相同维度 d' 的实验历史记录向量 $\hat{e}_i = (\hat{e}_{11}, \hat{e}_{12}, \dots, \hat{e}_{d'})$.

为了要让 \hat{e}_i 的维度匹配真实复杂系统的输入变量的维度, 还需要对 \hat{e}_i 做样本恢复和去除 \hat{y}_i . 样本恢复的步骤如下.

(1) 判断样本向量中各个部分原先属于哪个属性并进行恢复;

1) 对样本向量中经过归一化的部分进行归一化恢复;

2) 对样本向量中经过独热编码的部分进行独热编码恢复;

3) 对样本向量中经过标签编码的部分进行标签编码恢复;

(2) 对步骤(1)中得到的结果进行向量组装, 得到经过恢复后的实验样本.

在 Python 里, 这些恢复函数可以用各自原函数的 `inverse_transform` 方法实现. 经过恢复后的实验历史记录形如 (\hat{x}_i, \hat{y}_i) , 但由于 \hat{y}_i 由变分自编码器生成, 并不是系统的真实响应, 实验历史记录中必须要去除这个 \hat{y}_i . 最终, 由实验设计模块生成的样本记为 \hat{x}_i .

依据实验设计需求, 如果实验需要做 k 次, 则需要

让实验设计模块的流程进行 k 次, 生成 k 个实验样本, 这 k 个实验样本组成的集合 $O_g = (\hat{x}_1, \hat{x}_2, \dots, \hat{x}_k)$ 即为本文方法生成的实验设计方案.

2 实验设计有效性评价方法

为评价本文实验设计方法的有效性, 可用生成的实验样本和它们对应的实验响应去拟合真实复杂系统, 并评估其误差. 由于支持向量回归 (support vector regression, SVR) 模型^[19]在表达式形式未知和小样本的情况下具有良好的性能, 这里选用 SVR 来进行拟合.

SVR 是支持向量机^[20]在回归问题上的应用, 将输入映射到高维空间, 并用一个或一组超平面来对样本点进行回归. 假定在高维空间中, 样本点是可以通过如下线性函数拟合的:

$$\hat{y} = \hat{f}(x) = w \cdot x + b \quad (9)$$

其中, w 和 b 是该线性函数的参数, SVR 的损失函数可以表达成:

$$L(y, \hat{f}(x)) = \|w\|^2/2 + C \sum_{i=1}^N (\xi_i^- + \xi_i^+) \quad (10)$$

其中, C 是可调的超参数, N 是输入样本的数量, ξ_i^- 和 ξ_i^+ 是松弛变量. SVR 的优化目标是通过优化参数 w 和 b 来最小化它的损失函数, 即:

$$w, b = \arg \min_{w, b} L(y, \hat{f}(x)) \quad (11)$$

SVR 虽然将输入样本映射到了高维空间,但向量内积在原空间上更容易被计算出来,因此在高维空间上的向量内积运算使用核函数技巧^[21].常见的核函数包括:线性核、多项式核、径向基核和 Sigmoid 核.

利用 SVR 对实验设计评价方法如下:

首先,在原始样本空间中采样 n 个样本点,对每一个样本点 \mathbf{x}_i 实施实验,得到真实响应为: y_i .将 \mathbf{x}_i 和 y_i 组合得到真实样本和真实响应组合的集合 $\mathbf{O}_t = ((\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n))$,用作测试 SVR 拟合效果的测试集.

然后,对每一个待测实验设计方法 E_j ,令其生成含有 k 个样本的实验设计方案 \mathbf{O}_{gE_j} ,对 \mathbf{O}_{gE_j} 中每一个样本点 $\hat{\mathbf{x}}_i$ 实施实验,得到真实响应 $f(\hat{\mathbf{x}}_i)$.对每一个 E_j ,将对应的 $\hat{\mathbf{x}}_i$ 和 $f(\hat{\mathbf{x}}_i)$ 组合而成的集合 $\mathbf{O}_{sE_j} = ((\hat{\mathbf{x}}_1, f(\hat{\mathbf{x}}_1)), (\hat{\mathbf{x}}_2, f(\hat{\mathbf{x}}_2)), \dots, (\hat{\mathbf{x}}_k, f(\hat{\mathbf{x}}_k)))$,用作 SVR 拟合复杂系统的训练集.

对每一个 \mathbf{O}_{sE_j} 用 SVR 拟合,得到的模拟复杂系统的函数记为: $y'_{E_j} = f'_{E_j}(\mathbf{x})$.

将测试集中每个 \mathbf{x}_i 带入每个模拟的复杂系统,得到每一个 \mathbf{x}_i 的预测响应 $y'_{E_{j_i}} = f'_{E_{j_i}}(\mathbf{x}_i)$.

最后比较每种模拟的复杂系统的对测试集中所有 \mathbf{x}_i 的预测响应 $y'_{E_{j_i}}$ 和真实响应 y_i 之间的 $RMSE$ 、 MRE 、 MAE 、 R^2 , $RMSE$ 、 MRE 、 MAE 越低, R^2 越接近于1,说明拟合出的模拟复杂系统更接近真实复杂系统,即其对应的实验设计方法表现更好.

3 实验分析

本节选取了几种常用的实验设计方法作为基线,然后介绍了所使用的复杂系统模型和实验设置.最后通过实验得出了结果并对其进行了分析.

3.1 基线

基线实验设计选择拉丁超立方、K-means^[22]和

$$y_Q = 20.1168 \times (D - \Delta D) \frac{\sin(\varphi - \Delta\varphi + \Delta\beta) - \frac{V - \Delta V}{V_T - \Delta V_T} \sin(X - \Delta\beta + \Delta C)}{\sin(X + \varphi - \Delta\varphi + \Delta C)} \quad (13)$$

3.3 实验设置

在本次实验中,变分自编码器的隐变量空间的维度 l 设定为256,输入层里维度扩大的倍数 m 设定为5000,变分自编码器训练集 \mathbf{O}_h 的大小设定为1000,用标准正态分布随机抽样对隐变量 z 采样,SVR的测试集 \mathbf{O}_t 的大小设定为100,SVR的核函数设为线性核函数.

Maximin^[23]实验设计方法与本文方法相比较,因为它们是常用的实验设计方法且支持在连续区间内的实验取点.上述的这些基线实验设计方法可以通过 Python 实验设计库 doepy (<https://github.com/tirthajyoti/doepy>) 实现.

3.2 复杂系统模型

直航鱼雷在现代战争中是反制敌方舰船和潜水艇的有力武器,为了提高直航鱼雷的命中率和增加直航鱼雷的打击范围,需要考虑多种影响直航鱼雷发射的因素,建立直航鱼雷的命中模型并通过大量的仿真实验,得到清晰的影响直航鱼雷的命中率和打击范围的关键因素和它们之间的关系,从而可以为实际作战中直航鱼雷的战术运用提供指导.

本文实验选用文献[24]中的直航鱼雷命中模型作为本文实验设计方法验证的复杂系统.直航鱼雷命中模型的响应 y_Q 是鱼雷命中位置与打击目标中心的相对距离(m).该响应由以下的输入变量决定:目标速率(节) V 、鱼雷速率(节) V_T 、目标距离(链) D 、目标舷角 X 、发射鱼雷的提前角 φ 、目标距离观测误差 ΔD 、鱼雷航向误差角 $\Delta\varphi$ 、目标方位角观测误差 $\Delta\beta$ 、目标速率的观测误差 ΔV 、鱼雷速率的观测误差 ΔV_T 、目标航向角的观测误差 ΔC .

在本文实验中, $V = 15$, $V_T = 50$, $D = 50$, $X \in [-180, 180]$, $\Delta D \in [-15, 15]$, $\Delta\varphi \in [-3, 3]$, $\Delta\beta \in [-3, 3]$, $\Delta V \in [-3, 3]$, $\Delta V_T \in [-3, 3]$, $\Delta C \in [-3, 3]$. $N(\mu, \sigma)$ 是均值 μ 标准差 σ 的正态分布,其中, X 服从 $N(70, 5)$, ΔD 服从 $N(0, 5)$, $\Delta\varphi$ 服从 $N(0, 1)$, $\Delta\beta$ 服从 $N(0, 1)$, ΔV 服从 $N(0, 1)$, ΔV_T 服从 $N(0, 1)$, ΔC 服从 $N(0, 1)$. φ 由输入变量 V 、 V_T 、 X 决定:

$$\varphi = \arcsin\left(\frac{V}{V_T} \sin X\right) \quad (12)$$

整个直航鱼雷命中模型的方程如式(13)所示:

对每一种待测实验设计方法,生成的实验设计方案 \mathbf{O}_{gE_j} 的样本数大小 k 分别取5、10、15、20、25、30、35、40、45、50.限定 k 的最大值为50,即最多生成50个样本.

3.4 实验结果

实验结果如图5所示.图5是不同实验设计方法

在不同样本数下分别按图 5(a) $RMSE$ 、图 5(b) MRE 、图 5(c) MAE 、图 5(d) R^2 决定系数这些指标对直航鱼雷命中模型拟合效果比较的折线图。

从实验结果可以看出,对直航鱼雷命中模型的拟合,不论是在小样本还是在一定样本数下,本文方法均优于 K-means 实验设计方法, Maximin 实验设计方法与拉丁超立方实验设计方法。只有当样本数 k 为 15 时, Maximin 实验设计方法的 $RMSE = 28.96$ 、 $MAE = 23.55$ 与本文方法在样本数 k 为 15 时的 $RMSE = 27.49$ 、 $MAE = 22.71$ 相差很小,拟合效果相近。随着样本数 k 的增加,本文方法的 $RMSE$ 、 MRE 和 MAE 均能稳步减少,当样本数 k 等于 50 时,本文方法的 $RMSE = 4.30$ 、 $MRE = 12.28$ 、 $MAE = 3.35$ 均达到了各指标里所有实验验证结果中的

最低值。其他基线实验设计方法随着样本数 k 的增加,虽然它们的 $RMSE$ 、 MRE 和 MAE 总体也是呈下降的趋势,但并不稳定,有可能会在局部产生不降反升的情况,这也意味着本文方法不同于其他基线实验设计方法,随着样本数 k 的增加并不会让 SVR 显著的产生过拟合。从 R^2 决定系数上看,在小样本情况下,本文方法的 R^2 决定系数稍稍大于 0,其他基线方法的 R^2 决定系数均小于 0,它们的预测响应均和真实响应的相关性很差,当样本数 k 增加时,本文方法和基线实验设计方法均能提高它们的 R^2 决定系数并使它们最终都可以大于 0,但只有本文方法的 R^2 决定系数最终可以达到 0.7 以上,即预测响应和真实响应高度相关,在样本数 k 为 50 时,本文方法的 R^2 决定系数更是高达 0.98。

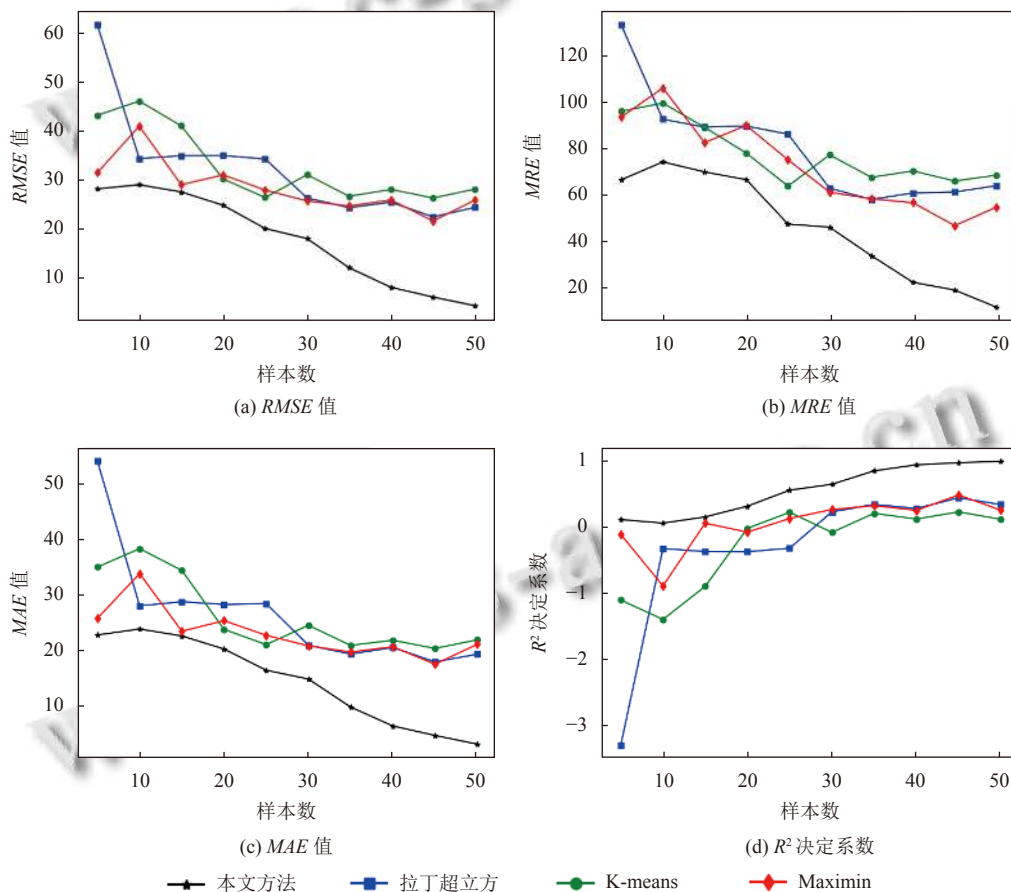


图 5 不同实验设计方法在不同指标不同样本数下对直航鱼雷命中模型拟合效果的比较

表 1 是各个实验设计方法的指标在低于或高于特定值时最少所需的样本数 k 。表 1 中的未达到表示在最多取 50 个样本的情况下,指标仍然未低于或高于所要求的特定的数值。由表 1 可知,本文方法总能以最少的

样本量获得最佳的实验结果,即说明本文的方法和其他基线方法相比是最高效的。

综上所述,本文实验设计方法可以提高在复杂系统中的实验设计效率,在拟合复杂系统时,本文实验设

计方法随着样本数的增加也不会产生显著的过拟合,即使在小样本的情况下,本文方法也能得到比较优秀的实验设计方案。

4 结论

本文为解决复杂系统实验设计中非线性、不确定性、高相关性和规模性导致的维度灾难问题,提出了一种基于变分自编码器的复杂系统实验设计方法。首先,给出了基于变分自编码器的复杂系统实验设计方法的总体结构,详细介绍了每个子模块的目的和它们的功能,并给出了评价实验设计优劣的方法。然后,本文通过对实际的直航鱼雷命中模型的拟合,对本文方法、拉丁超立方、K-means 和 Maximin 实验设计方法进行了比较实验。最后,得出了本文方法可以提高复杂系统实验效率的结论。

表1 各个实验设计方法的指标在低于或高于特定值时最少所需的样本数 k

指标	条件	本文方法	拉丁超立方	K-means	Maximin
RMSE	<50	5	10	5	5
	<40	5	10	20	5
	<30	5	30	25	15
	<20	30	未达到	未达到	未达到
	<10	40	未达到	未达到	未达到
MRE	<100	5	10	5	5
	<80	5	30	20	25
	<60	25	35	未达到	35
	<40	35	未达到	未达到	未达到
MAE	<40	5	10	5	5
	<30	5	10	20	5
	<20	25	35	未达到	35
	<10	35	未达到	未达到	未达到
R^2 决定系数	>0	5	30	25	15
	>0.15	15	30	25	30
	>0.30	20	35	未达到	35
	>0.45	25	未达到	未达到	45
	>0.60	30	未达到	未达到	未达到
	>0.75	35	未达到	未达到	未达到

参考文献

- 1 万秋生. 复杂仿真实验设计与监控技术研究 [硕士学位论文]. 哈尔滨: 哈尔滨工业大学, 2012. [doi: 10.7666/d.D241030]
- 2 Hsu MC, Akkerman I, Bazilevs Y. Finite element simulation of wind turbine aerodynamics: Validation study using NREL Phase VI experiment. *Wind Energy*, 2014, 17(3): 461–481. [doi: 10.1002/we.1599]

- 3 卞曰塘, 刘夏群, 李金生. 基于网络混合学习策略的股市投资者行为演化模型及仿真. *运筹与管理*, 2019, 28(11): 156–168.
- 4 郭啸天, 李文波, 王海雷, 等. 基于 ABM 与 ABS 的暴雨洪涝人口风险动态模拟. *计算机系统应用*, 2015, 24(12): 10–17. [doi: 10.3969/j.issn.1003-3254.2015.12.002]
- 5 Jomon G, Jojish JV, Santhanakrishnan T. System of systems architecture for generic torpedo defence system for surface ships. *Advances in Military Technology*, 2019, 14(2): 307–319. [doi: 10.3849/aimt.01330]
- 6 谢晓天. 复杂仿真系统高效实验方法研究 [硕士学位论文]. 哈尔滨: 哈尔滨工业大学, 2016. [doi: 10.7666/d.D01098648]
- 7 Shokri A. Application of Sono-photo-Fenton process for degradation of phenol derivatives in petrochemical wastewater using full factorial design of experiment. *International Journal of Industrial Chemistry*, 2018, 9(4): 295–303. [doi: 10.1007/s40090-018-0159-y]
- 8 Yolmeh M, Jafari SM. Applications of response surface methodology in the food industry processes. *Food and Bioprocess Technology*, 2017, 10(3): 413–433. [doi: 10.1007/s11947-016-1855-2]
- 9 Bessa MA, Bostanabad R, Liu Z, *et al.* A framework for data-driven analysis of materials under uncertainty: Countering the curse of dimensionality. *Computer Methods in Applied Mechanics and Engineering*, 2017, 320: 633–667. [doi: 10.1016/j.cma.2017.03.037]
- 10 Feng GZ, Lei SY, Guo YJ, *et al.* Optimisation of air-distributor channel structural parameters based on Taguchi orthogonal design. *Case Studies in Thermal Engineering*, 2020, 21: 100685. [doi: 10.1016/j.csite.2020.100685]
- 11 Feng ZK, Niu WJ, Cheng CT. Optimizing electrical power production of hydropower system by uniform progressive optimality algorithm based on two-stage search mechanism and uniform design. *Journal of Cleaner Production*, 2018, 190: 432–442. [doi: 10.1016/j.jclepro.2018.04.134]
- 12 Ma P, Zhou YC, Shang XB, *et al.* Firing accuracy evaluation of electromagnetic railgun based on multicriteria optimal Latin hypercube design. *IEEE Transactions on Plasma Science*, 2017, 45(7): 1503–1511. [doi: 10.1109/TPS.2017.2705980]
- 13 Kingma DP, Welling M. Auto-encoding variational Bayes. arXiv: 1312.6114, 2014.
- 14 Hou XX, Shen LL, Sun K, *et al.* Deep feature consistent variational autoencoder. *Proceedings of 2017 IEEE Winter Conference on Applications of Computer Vision*. Santa

- Rosa: IEEE, 2017. 1133–1141. [doi: [10.1109/WACV.2017.131](https://doi.org/10.1109/WACV.2017.131)]
- 15 Hsu WN, Zhang Y, Glass J. Unsupervised domain adaptation for robust speech recognition via variational autoencoder-based data augmentation. Proceedings of 2017 IEEE Automatic Speech Recognition and Understanding Workshop. Okinawa: IEEE, 2017. 16–23. [doi: [10.1109/ASRU.2017.8268911](https://doi.org/10.1109/ASRU.2017.8268911)]
- 16 Semeniuta S, Severyn A, Barth E. A hybrid convolutional variational autoencoder for text generation. Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing. Copenhagen: Association for Computational Linguistics, 2017. 627–637.
- 17 Pedregosa F, Varoquaux G, Gramfort A, *et al.* Scikit-learn: Machine learning in Python. The Journal of Machine Learning Research, 2011, 12: 2825–2830.
- 18 翟正利, 梁振明, 周炜, 等. 变分自编码器模型综述. 计算机工程与应用, 2019, 55(3): 1–9. [doi: [10.3778/j.issn.1002-8331.1810-0284](https://doi.org/10.3778/j.issn.1002-8331.1810-0284)]
- 19 Mohammadi B, Mehdizadeh S. Modeling daily reference evapotranspiration *via* a novel approach based on support vector regression coupled with whale optimization algorithm. Agricultural Water Management, 2020, 237: 106145. [doi: [10.1016/j.agwat.2020.106145](https://doi.org/10.1016/j.agwat.2020.106145)]
- 20 林香亮, 袁瑞, 孙玉秋, 等. 支持向量机的基本理论和研究进展. 长江大学学报(自科版), 2018, 15(17): 48–53.
- 21 肖建, 于龙, 白裔峰. 支持向量回归中核函数和超参数选择方法综述. 西南交通大学学报, 2008, 43(3): 297–303. [doi: [10.3969/j.issn.0258-2724.2008.03.001](https://doi.org/10.3969/j.issn.0258-2724.2008.03.001)]
- 22 Shakeel PM, Baskar S, Dhulipala VRS, *et al.* Cloud based framework for diagnosis of diabetes mellitus using K-means clustering. Health Information Science and Systems, 2018, 6(1): 16. [doi: [10.1007/s13755-018-0054-0](https://doi.org/10.1007/s13755-018-0054-0)]
- 23 Pronzato L. Minimax and maximin space-filling designs: Some properties and methods for construction. Journal de la Société Française de Statistique, 2017, 158(1): 7–36.
- 24 孙华春, 李长文, 李海玲. 直航鱼雷命中概率模型与仿真. 舰船电子工程, 2009, 29(12): 138–141. [doi: [10.3969/j.issn.1627-9730.2009.12.038](https://doi.org/10.3969/j.issn.1627-9730.2009.12.038)]