

Coverage-Based Designs Improve Sample Mining and Hyper-Parameter Optimization

Gowtham Muniraju[†], Bhavya Kailkhura[‡], Jayaraman J. Thiagarajan[‡], Peer-Timo Bremer[‡], Cihan Tepedelenlioglu[†], *Senior Member, IEEE* and Andreas Spanias[†], *Fellow IEEE*.

Abstract—Sampling one or more effective solutions from large search spaces is a recurring idea in machine learning, and sequential optimization has become a popular solution. Typical examples include data summarization, sample mining for predictive modeling and hyper-parameter optimization. Existing solutions attempt to adaptively trade-off between global exploration and local exploitation, wherein the initial exploratory sample is critical to their success. While discrepancy-based samples have become the *de facto* approach for exploration, results from computer graphics suggest that coverage-based designs, e.g. Poisson disk sampling, can be a superior alternative. In order to successfully adopt coverage-based sample designs to ML applications, which were originally developed for $2 - d$ image analysis, we propose fundamental advances by constructing a parameterized family of designs with provably improved coverage characteristics, and by developing algorithms for effective sample synthesis. Using experiments in sample mining and hyper-parameter optimization for supervised learning, we show that our approach consistently outperforms existing exploratory sampling methods in both blind exploration, and sequential search with Bayesian optimization.

Index Terms—Hyper-parameter optimization, coverage-based sample design, Poisson disk sampling, predictive modeling, sequential optimization.

I. INTRODUCTION

A. Sampling in Machine Learning

Sample design has been a long-standing research area in statistics [1], and has now become a crucial problem in machine learning and AI, particularly with the emergence of numerous data-driven learning paradigms. The notion of sampling appears in a variety of contexts, ranging from summarizing complex data [2], generating mini-batches for effective neural network training [3], metric learning [4] to hyper-parameter search [5], [6], reinforcement learning [7], [8] and knowledge transfer [9]. A common goal in these seemingly diverse applications is to identify one or more effective solutions from a large search space, using the smallest amount of resources. In principle, there are two competing strategies while performing sampling [10]: *exploitation*, which probes a limited region in the search space with the hope of improving an already identified solution; and *exploration*, which probes a larger part of the search space with the hope of finding solutions that are yet to be refined. In practice, sequential sampling methods that can trade-off between exploration and exploitation strategies are preferred [11]. However,

given the large volume of typical search spaces and restrictions on resources (time and compute), the exploration step is highly critical to reduce the uncertainties to an extent that the exploitation step can be expected to succeed. Over the last several decades, a large class of exploratory sampling techniques have been developed [12]–[14]. Though the overarching objective is to cover the search space uniformly, it is well known that uniformity alone does not suffice. For example, optimal sphere packings lead to highly uniform designs, yet are prone to causing aliasing artifacts. Consequently, effective exploration requires to balance uniformity and randomness in the search space, often evaluated using heuristic measures such as discrepancy [15]. More recently, the *pair correlation function* (PCF) has been found to be a more useful statistic for evaluating the quality of sample designs [16]–[19].

While discrepancy-based quasi-random designs have been commonly utilized in several applications [5], [20], the computer graphics community has had long-standing success with coverage-based designs, in particular Poisson Disk Sampling (PDS) [21], [22]. The works in [21], [22] were the first to introduce PDS for turning regular aliasing patterns into featureless noise, which makes them perceptually less visible. Their works were inspired by the seminal work of Yellott et.al. [23], who observed that the photo-receptors in the retina of monkeys and humans are distributed according to a Poisson disk distribution. For the first time in [17], PDS was formally defined using the pair correlation function and used to obtain theoretical bounds on achievable coverage properties. Despite their well-established success in image/volume rendering [12], [24], [25], coverage-based designs have not been adopted in the machine learning community. Recently, in [19], Kailkhura *et al.* developed a generic spectral sampling framework, that encompasses several existing designs including PDS, blue noise [25] and variants [18], by jointly analyzing the spatial and spectral properties of sample distributions. Though this framework enjoys several desirable properties in theory, constructing an optimal design and actually synthesizing samples that match these characteristics are challenging. When designed sub-optimally, a spectral sample can perform worse than other random sampling strategies. In addition, the sample synthesis is based solely on PCF matching, which is a summary 1–D statistic of high-dimensional point clouds, thus making this optimization very challenging in practice.

B. Proposed Work

In this work, we propose to develop novel coverage-based designs for challenging machine learning problems,

[†] G. Muniraju, C. Tepedelenlioglu and A. Spanias are with School of ECEE, Arizona State University. Email: {gmuniraj, cihan, spanias}@asu.edu.

[‡] B. Kailkhura, J. J. Thiagarajan and P-T. Bremer are with the Center for Applied Scientific Computing, Lawrence Livermore National Laboratory. Email: {kailkhura1, jjayaram, bremer5}@llnl.gov

namely sample mining in predictive modeling and hyper-parameter optimization. Building upon the theoretical foundations from [19], we argue that larger coverage is critical to improving the expected performance of sample designs in ML tasks that require the exploration of complex optimization surfaces. Further, we make the following key contributions to produce highly effective samples in practice:

- We introduce a new parameterized family of coverage-based designs using the pair correlation function, which generalizes existing constructions such as [17], [19], for machine learning applications;
- Using tools from spectral sampling theory, we show that the proposed sample design achieves the largest coverage so far;
- For the first time, we develop an efficient strategy to find the “optimal” parameters of a design (i.e. with largest coverage) for a given sample size and dimensionality;
- We design a scalable and effective sample synthesis algorithm that consistently outperforms existing PCF matching approaches such as [16] and [19].
- Using empirical studies on predictive modeling, we demonstrate that the proposed sample design outperforms existing discrepancy-based and coverage-based designs, particularly under low sampling rates.

The proposed coverage-based design is based on systematically trading-off randomness characteristics of a point distribution with coverage to enable improved performance. Such a controlled random sampling is mathematically represented using the PCF and analyzed via the spectral sampling principles from [19]. Surprisingly, we find that the achievable coverage of the proposed design is significantly larger than a conventional PDS design ($\sim 25\% - 40\%$ increase for the same configuration). Further, our synthesis algorithm consistently produces high-quality samples and is highly robust, as evidenced by the performance variance across multiple realizations.

In order to demonstrate the importance of coverage-based designs in challenging applications, we consider the problem of hyper-parameter tuning while building ML models. To this end, we consider scenarios where we rely solely on exploration (*blind sampling*), similar to [5], and where we use the exploratory samples to initialize a Bayesian optimization pipeline with expected improvement as the acquisition function, as carried out in [20]. We perform empirical studies with (a) a standard feature extractor-classifier pipeline, and (b) deep neural networks that perform end-to-end learning. Our results show that the proposed sample design consistently outperforms state-of-the-art exploratory sampling methods including Latin Hyper Cube (LHS), Quasi-Monte Carlo (QMC) designs [15] and spectral samples in [19]. Interestingly, we observed significant improvements even in the Bayesian optimization cases, which clearly emphasizes the importance of the initial exploration step. In summary, the effectiveness of our approach even with small sample sizes establishes it as a powerful exploratory sampling technique for ML/AI applications.

II. COVERAGE-BASED SAMPLE DESIGNS

Though a variety of discrepancy measures are commonly used for exploratory sampling, our focus is on coverage-based

designs. In this section, we briefly describe the mathematical tools required for the design and analysis of coverage-based sampling. Subsequently, we discuss two popular coverage-based designs from [17] and [19] respectively.

Broadly speaking, a reasonable objective for exploratory sampling is to ensure that the samples are random, thus providing an equal chance of finding meaningful solutions anywhere in the search space. However, in order to ensure diversity, a second objective is often considered, which is to cover the space uniformly. In this paper, we consider the general class of coverage-based sample designs [19]:

Definition 1. (Coverage-based Design) A set of N random samples $\{\mathbf{X}_i\}_{i=1}^N$ in a search space \mathcal{D} can be characterized as a coverage-based design, if $\{\mathbf{X}_i = \mathbf{x}_i \in \mathcal{D}; i = 1, \dots, N\}$ satisfy the following two objectives:

- $\forall \mathbf{X}_i, \forall \Delta \mathcal{D} \subseteq \mathcal{D} : Pr(\mathbf{X}_i = \mathbf{x}_i \in \Delta \mathcal{D}) = \frac{1}{\Delta \mathcal{D}} \int_{\Delta \mathcal{D}} d\mathbf{x};$
- $\forall \mathbf{x}_i, \mathbf{x}_j : \|\mathbf{x}_i - \mathbf{x}_j\| \geq r_{\min},$

where r_{\min} is referred to as the *coverage radius* (or disk size). In this definition, the first objective states that the probability of a random sample \mathbf{X}_i falling inside a subset $\Delta \mathcal{D}$ of \mathcal{D} is equal to the hyper-volume of $\Delta \mathcal{D}$. The second condition enforces the disk constraint for improving coverage. Since, there existed no quality metrics to jointly characterize the coverage and randomness properties, several recent works have adopted the pair correlation function [16] as a quality metric.

Definition 2. (Pair Correlation Function) Let us denote the intensity of a point process \mathcal{X} as $\lambda(\mathcal{X})$, i.e., the average number of points in an infinitesimal volume around \mathcal{X} . For isotropic point processes, this is a constant. To define the product density β , let $\{B_i\}$ denote the set of infinitesimal spheres around the points, and $\{dV_i\}$ indicate the volume measures of B_i . Then, we have $Pr(\mathbf{X}_1 = \mathbf{x}_1, \dots, \mathbf{X}_N = \mathbf{x}_N) = \beta(\mathbf{x}_1, \dots, \mathbf{x}_N) dV_1 \dots dV_N$ which represents the probability of having points \mathbf{x}_i in $\{B_i\}$. In the isotropic case, for a pair of points, β depends only on the distance between them, and hence $\beta(\mathbf{x}_i, \mathbf{x}_j) = \beta(\|\mathbf{x}_i - \mathbf{x}_j\|) = \beta(r)$ and $Pr(r) = \beta(r) dV_i dV_j$. The PCF is then defined as $G(r) = \beta/\lambda^2$.

Alternatively, Fourier analysis can be utilized for understanding the qualitative properties of sampling patterns. For isotropic samples, a metric of interest is the radially-averaged power spectral density (PSD), which describes how the signal power is distributed over spatial frequencies.

Definition 3. (Radially-averaged PSD) For a finite set of N points, $\{\mathbf{x}_j\}_{j=1}^N$, in a region with unit volume, the PSD of the sampling function $\sum_{j=1}^N \delta(\mathbf{x} - \mathbf{x}_j)$ is defined as

$$P(\mathbf{k}) = \frac{1}{N} |S(\mathbf{k})|^2 = \frac{1}{N} \sum_{j,\ell} e^{-2\pi i \mathbf{k} \cdot (\mathbf{x}_\ell - \mathbf{x}_j)}, \quad (1)$$

where $|\cdot|$ denotes the ℓ_2 -norm and $S(\mathbf{k})$ denotes the Fourier transform of the sampling function.

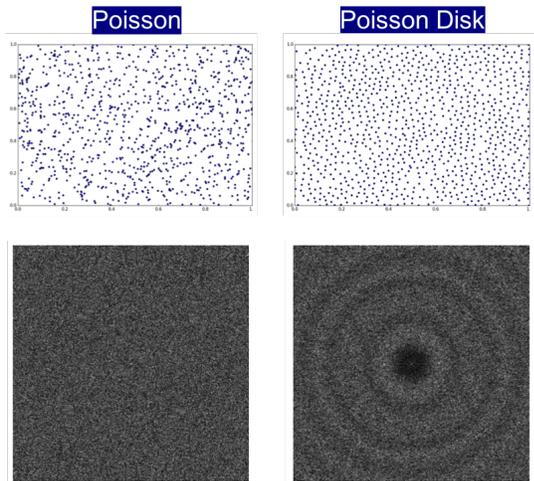


Figure 1: Illustration of 2-d patterns obtained using Poisson and Poisson disk sampling. We show the point distribution (top) and the power spectral density (bottom) for each case.

Interestingly, there is a well-defined connection between the PCF of a sample design and its radially-averaged PSD, and this connection is central to the proposed work.

Definition 4. (Linking PCF and PSD) For an isotropic sample design with N points, $\{\mathbf{x}_j\}_{j=1}^N$, in a d -dimensional region, the radially averaged power spectral density $P(k)$ and the pair correlation function $G(r)$ are related as follows:

$$P(k) = 1 + \frac{N}{V} (2\pi)^{\frac{d}{2}} k^{1-\frac{d}{2}} H_{\frac{d}{2}-1} \left[r^{\frac{d}{2}-1} G(r) - 1 \right], \quad (2)$$

where k is the frequency index, V is the volume of the sampling region and $H_d[\cdot]$ denotes the Hankel transform,

$$H_d(f(r))(k) = \int_0^\infty r J_d(kr) f(r) dr,$$

with $J_d(\cdot)$ denoting the Bessel function of order d .

Finally, it is important to note that, not every PCF construction is physically realizable by a sample design. In fact, there are two necessary mathematical conditions¹ that a design must satisfy to be realizable.

Definition 5. (Realizability) A PCF can be considered to be potentially realizable through a sample design, if it satisfies:

- the PCF must be non-negative, i.e., $G(r) \geq 0$, $\forall r$, and
- the corresponding PSD must be non-negative, i.e., $P(k) \geq 0$, $\forall k$.

A. Poisson Disk Sampling

The well-known Poisson design (Figure 1(a)) enforces only the first condition from Definition 1, in which case the number of samples that fall inside any subset $\Delta\mathcal{D} \subseteq \mathcal{D}$ obeys a discrete Poisson distribution. Consequently, Poisson disk sampling [17] (Figure 1(b)) that explicitly enforces the disk constraint is considered to be optimal in this context. Several

¹Whether or not these two conditions are not only necessary but also sufficient is still an open question (however, no counterexamples are known).

widely adopted strategies for generating Poisson disk samples rely on the heuristic idea of dart throwing [12], [13], [21], [22], which uses as many darts as required to cover the search space, while not violating the disk criterion. Despite its effectiveness, its primary shortcoming is the choice of termination condition, since it is not easy to quantify the coverage and randomness properties. This motivated the use of pair correlation function (PCF) [16] to summarize spatial characteristics of a sampling pattern, using which Kailkhura *et al.* [17] formally defined Poisson disk sampling for the first time (Figure 2(a)).

Definition 6. (Poisson disk sampling) [17] Given the desired disk size r_{\min} , PDS is defined using the PCF as

$$G(r - r_{\min}) = \begin{cases} 0 & \text{if } r < r_{\min} \\ 1 & \text{if } r \geq r_{\min}. \end{cases} \quad (3)$$

Note that, the disk radius r_{\min} is referred as the coverage.

B. Space-Filling Spectral Design

While Poisson disk sampling was preferred in computer graphics applications for turning regular aliasing patterns into featureless noise, it is not directly suitable for conventional predictive modeling problems. Consequently, in [19], the authors studied the impact of coverage on supervised regression problems, and provided empirical evidence that larger coverage in the sample design was critical to improving the expected performance of the models designed using them. Motivated by this observation, they developed the *space-filling spectral design* (SFSD), which is defined as:

Definition 7. (Space-Filling Spectral Design) [19]

$$G(r; r_{\min}, r_1, P_0) = f(r - r_1) + P_0 (f(r - r_{\min}) - f(r - r_1)), \quad (4)$$

$$\text{with } f(r - r_{\min}) = \begin{cases} 0 & \text{if } r \leq r_{\min} \\ 1 & \text{if } r > r_{\min} \end{cases},$$

where $r_{\min} \leq r_1$ and $P_0 \geq 1$.

This represents a parameterized stair function (Figure 2(b)) that introduces a peak in the quest of increasing the coverage.

III. PROPOSED SAMPLE DESIGN METHODOLOGY

In order to enable the effective use of coverage-based designs in ML applications, we need to confront the following challenges: (i) most existing methods are designed specifically for 2-d, and a trivial extension of such constructions provide poor coverage, even in $d > 3$; (ii) current sample synthesis algorithms based on PCF matching require extensive manual tuning, and perform poorly as the dimension increases – in many cases, the synthesis quality is no better than random sampling; and (iii) the superior performance of coverage-based designs has been established mostly on graphics tasks, such as image/volume rendering, and similar gains are yet to be achieved in ML applications.

In this section, we first propose a new parameterized PCF construction for coverage-based designs, which achieves larger

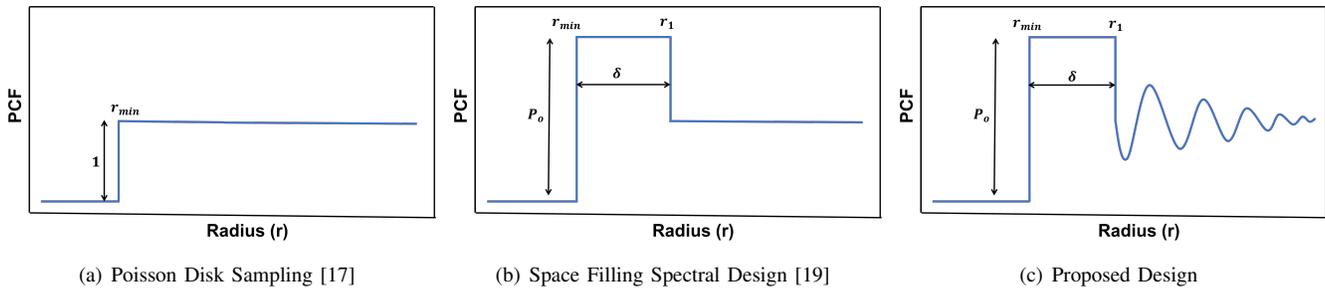


Figure 2: Pair correlation functions of coverage-based sample designs. Each design leads to a different coverage size r_{min} .

coverage compared to existing approaches. Next, we develop a practical strategy for finding an “optimal” PCF configuration. Finally, in the next section, we present an effective sample synthesis algorithm for coverage-based sampling, that consistently leads to high-quality samples across different sample sizes and dimensionality.

A. A New Parameterized Family

Following notations in the previous section, our parameterized PCF construction for coverage-based sampling can be expressed as follows:

$$G(r; r_{min}, r_1, P_0, A, B, C, D) = P_0 (f(r - r_{min}) - f(r - r_1)) + \left(1 + \frac{A}{r} \exp(-Br) \sin(2\pi Cr + D)\right) * f(r - r_1), \quad (5)$$

$$\text{where } f(r - r_{min}) = \begin{cases} 0 & \text{if } r \leq r_{min} \\ 1 & \text{if } r > r_{min} \end{cases},$$

$$r_{min} \leq r_1 \text{ and } P_0 \geq 1.$$

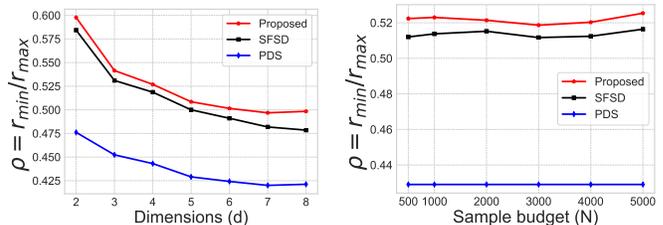
The intuition behind this construction is to enable trade-off between randomness and uniformity/coverage properties of a sample design. This construction (see Figure 2(c)) has three crucial properties:

- 1) The PCF is zero from $0 \leq r \leq r_{min}$, corresponding to the coverage size similar to other designs;
- 2) The PCF has a peak from $r_{min} < r \leq r_1$ and damped oscillations from $r > r_1$ characterizing randomness;
- 3) The peak height P_0 , width $\delta = r_1 - r_{min}$, and oscillations can be adjusted to control the randomness property of a design, which in turn can maximize the coverage r_{min} .

The radially-averaged power spectral density of the PCF in (5) can be obtained using the relation in (2). As we will show later, this connection is central for designing and optimizing the proposed design in a computationally efficient manner.

B. Quantifying Coverage Gain

Next, we evaluate the coverage gain in our proposed design with respect to other coverage-based approaches from the computer graphics and surrogate modeling literature. For this analysis, we varied the parameter P_0 in the range $[1.00, 2.5]$ and performed a brute-force search on the parameter $r_1 \in [r_{min}, 2r_{min}]$, $A \in [0.1, 0.9]$, $B \in [2, 6]$, $C \in [50, 600]$, $D \in [-\pi, \pi]$, such that r_{min} is maximized, while also ensuring that the realizability conditions from Definition 5 are met.



(a) Fixed $N = 1000$ and varying d (b) Fixed $d = 5$ and varying N

Figure 3: Maximum achievable *relative radius* using different coverage-based designs. The proposed design consistently outperforms PDS and SFSD approaches.

We compare the coverage characteristics of the proposed approach to existing coverage-based designs, namely PDS and SFSD. More specifically, we compare different coverage-based designs using the *relative radius* [26] $\rho = r_{min}/r_{max}$, where r_{max} is the maximum possible radius for N samples in d dimensions. For a given N and d , r_{max} can be computed as

$$r_{max} = \sqrt[d]{\gamma_d \frac{V \Gamma(\frac{d}{2} + 1)}{\pi^{\frac{d}{2}} N}},$$

where the maximum packing density γ_d for $d = \{2, \dots, 8\}$ can be found in [27].

In Figure 3(a), we first fixed the the sample size at $N = 1000$. Subsequently, we measured the maximum achievable relative radius using different coverage-based designs of size N in dimensions $d = \{2, \dots, 8\}$ respectively. The first striking observation is that by incorporating controlled randomness, both SFSD and the proposed design produce significantly larger relative radius when compared to the conventional PDS. Further, in all cases, the proposed approach provides improved coverage over SFSD and as we will demonstrate in our results, this seemingly marginal improvement leads to significant performance gains in practice. Another interesting observation is that as the dimensionality increases, the relative radius decreases rapidly and all coverage-based design behave similarly. In other words, due to the curse of dimensionality, when the search space is comprised of tens of dimensions, the proposed approach will become similar to PDS (unless the sample size grows exponentially), while still being superior to discrepancy-based designs.

Algorithm 1 Automatic selection of PDS parameters

```

1: Input: Number of samples  $N$ , dimension  $d$ , parameter  $P_0$ ,
   step size  $\lambda$ ,  $V = 1$ .
2:  $\bar{r}_{min} = \sqrt[d]{V\Gamma(d/2 + 1)/(\pi^{d/2}N)}$   $\triangleright$  Conventional PDS
3: Initialize:  $r_{min} = \bar{r}_{min}, r_1 = 2r_{min}$ 
4:  $G(r) \leftarrow G(r; r_{min}, r_1, P_0)$   $\triangleright$  Initialize PCF
5:  $P(k) \leftarrow 1 + \frac{N}{V}(2\pi)^{\frac{d}{2}}k^{1-\frac{d}{2}}H_{\frac{d}{2}-1}\left[r^{\frac{d}{2}-1}G(r) - 1\right]$   $\triangleright$ 
   from (2)
6:  $k^* \leftarrow \arg \min_k P(k)$ 
7: While ( $P(k^*) \geq 0$ )  $\triangleright$  Constraint for realizable PCF
8:   Update  $r_{min}$ 
9:    $r_{min} \leftarrow r_{min} + \lambda\left(r_{min}\frac{\partial}{\partial r_{min}}P(k^*) + P(k^*)\right)$ 
10:  Update  $r_1$ 
11:   $r_1 \leftarrow r_1 - \lambda\left(r_{min}\frac{\partial}{\partial r_1}P(k^*)\right)$ 
12: Return  $r_{min}, r_1$   $\triangleright$  Optimal PCF settings

```

In practice, with applications such as hyper-parameter optimization, unless the intrinsic dimension of the optimization surface is low, exploratory sampling will be ineffective (even with million samples) as the volume of spaces grows exponentially with dimension. In the literature, it has been observed that the intrinsic dimension of search spaces is often between 3 – 6 over different datasets [5]. Similarly, in Figure 3(b), we show the relative radius ρ at different sample sizes for a given dimension $d = 5$. As it can be seen, for all coverage-based designs, the best achievable relative radius is nearly a constant at all sample sizes. Furthermore, the proposed design consistently produces larger coverage compared to other designs in all cases.

C. A Practical Strategy for Parameter Selection

A typical approach to find parameters that achieve the largest coverage gain is a brute-force search [19]. However, the search space of realizable parameters is complex – non-monotonic, coupled with the need to satisfy realizability conditions. To overcome this challenge, we develop an efficient gradient based parameter selection strategy for optimal PCF construction. Specifically, we are interested in solving the following parameter search problem²:

$$\begin{aligned}
& \text{maximize} : r_{min} \\
& \text{subject to} : P(k) \geq 0, \forall k \\
& \quad r_1 > r_{min}
\end{aligned} \tag{6}$$

Since the goal is to achieve maximal coverage, we maximize r_{min} such that the resulting PCF is realizable, which is verified by ensuring that the power spectral density $P(k) \geq 0, \forall k$. In our experiments, we found that the lagrangian relaxation of (6) is hard to optimize. Instead, we maximize an alternative objective function $r_{min} \times P(k^*)$ where $k^* = \arg \min_k P(k)$,

²We found that the maximum achievable coverage, r_{min} , for a given d and N , depends primarily on the choice of r_1 and P_0 , while the choices for other parameters A, B, C, D are not particularly sensitive. Thus, we only optimize over r_1 and P_0 in this paper.

Algorithm 2 Sample Design using GD-ALR Algorithm

```

1: Input: Number of samples  $N$ , dimension  $d$ , target PCF
    $\hat{G}^*(r_j)$ , learning rate  $\lambda$ 
2:  $\mathbf{X} \leftarrow \text{Random}(N, d)$   $\triangleright$  Initial random sample design
3:  $G \leftarrow \text{PCF}(\mathbf{X})$   $\triangleright$  Calculate initial PCF using Eq. (8)
4: for  $t = 1$  to  $T$  do  $\triangleright$  Total  $T$  gradient descent iterations
5:   for  $i = 1$  to  $N$  do  $\triangleright$  Update each sample at a time
6:      $\Delta_i^p \leftarrow \frac{\partial}{\partial x_i^p} \sum_{j=1}^M (G^t(r_j) - G^*(r_j))^2$  for  $p \in$ 
        $\{1, \dots, d\}$   $\triangleright$  Calculate gradients
7:      $\lambda \leftarrow 0.1e^{-0.1\sqrt{t}}$   $\triangleright$  Adapting learning rate
8:      $\mathbf{x}_i^p(t+1) \leftarrow \mathbf{x}_i^p(t) - \lambda \frac{\Delta_i^p}{|\Delta_i^p|}$   $\triangleright$  Update the samples
       position
9:      $G^t \leftarrow \text{PCF}(\mathbf{X})$   $\triangleright$  Update the PCF
10: return  $\mathbf{X}$   $\triangleright$  Optimized Samples

```

(i.e., consider only the minimum value of $P(k)$) is found to work better. In Algorithm 1, we present our approach to solve this modified optimization problem.

IV. PROPOSED SYNTHESIS ALGORITHM

We develop an approach that iteratively transforms an initial random sample design such that its PCF matches the PCF of the optimal coverage-based design. More specifically, we consider a non-linear least squares formulation similar to [16], [19]. Despite being computationally efficient, due to the high non-convexity of the PCF matching problem, conventional gradient descent based approaches perform poorly as the dimension increases. In fact, due to the small effective r_{min} , the synthesis quality is no better than random sampling. Here, we adopt a different approach to alleviate this limitation, and make PCF matching-based synthesis a viable solution.

Denoting the desired PCF for an optimal design by $G^*(r)$, we discretize the radius r into M points $\{r_j\}_{j=1}^M$, and minimize the sum of the weighted squares of errors between the target PCF $G^*(r_j)$ and the curve-fit function (explained next) $G(r_j)$. Consequently, sample synthesis is posed as the following non-linear least squares problem:

$$\min \sum_{j=1}^M (G(r_j) - G^*(r_j))^2. \tag{7}$$

A. PCF Matching Algorithm

Intuitively, the proposed PCF matching algorithm is comprised of two phases: achieving coverage and matching oscillations. In the first phase, the initial design of uniform random samples is optimized to achieve coverage by shifting the positions of the N samples, such that no two samples are closer than r_{min} . In the second phase, samples are optimized to match oscillations in the target PCF. Before presenting the proposed algorithm, we first describe the PCF estimator employed in our optimization.

PCF Estimator: To estimate the PCF of point samples, we employ a kernel density estimator [16], defined as

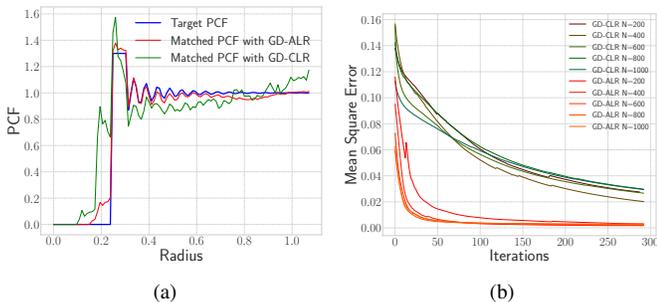


Figure 4: Coverage-based Sample Synthesis. (a) PCF matching performance of GD-ALR versus GD-CLR, while obtaining maximal coverage (r_{min}) for $N = 200$, $d = 4$, $P_0 = 1.3$. (b) Mean square error of PCF matching obtained using GD-ALR and GD-CLR, across different gradient descent iterations, for a fixed dimension $d = 4$.

$$\hat{G}(r) = \frac{V_W V_W}{\gamma_W N S_E(N-1)} \frac{1}{\sum_{i=1}^N \sum_{\substack{j=1 \\ i \neq j}}^N} k(r - |\mathbf{x}_i - \mathbf{x}_j|) \quad (8)$$

where $k(\cdot)$ denotes the Gaussian kernel function, $k(z) = (1/\sqrt{\pi}\sigma) \exp(-z^2/2\sigma^2)$. In this expression, V_W indicates the volume of the search space and S_E denotes the area of hypersphere. Finally, γ_W is an isotropic set covariance function which can be approximated as $\gamma_W = V_W - (S_W/\pi)r$, where S_W denotes the surface area of the sampling region. The term $\frac{V_W}{\gamma_W}$ accounts for edge correction for the unboundedness of the estimator.

Algorithm: Given the PCF estimate, the matching problem can be solved using gradient descent. However, due to the highly non-convex nature of this problem, gradient descent with constant learning rate (GD-CLR) [19] perform very poorly. Instead, we propose to employ gradient descent with adaptive learning rate, GD-ALR (Algorithm 2), with the learning rate update rule: $\lambda = 0.1e^{-0.1\sqrt{t}}$, for iteration t . More importantly, we find that, in order to achieve the maximal coverage, it is important to first optimize for coverage (updates with larger values of λ) and, then for oscillations (updates with smaller value of λ), instead of joint optimization as done in existing approaches [19], [25]. This behavior is illustrated in Fig. 4(a). Separately optimizing for coverage/oscillations using an adaptive learning rate profile, solves a major bottleneck in synthesizing coverage-based designs. In particular, we found that many other variants of gradient descent (e.g. Levenberg-Marquardt) failed to achieve the desired performance.

From Fig. 4(b), it can be seen that the proposed GD-ALR demonstrates superior convergence characteristics when compared to GD-CLR. We also conducted experiments with other optimization approaches, such as, momentum gradient descent optimizer. We observed that in all settings of N and d , GD-ALR outperformed other optimizers with faster convergence and significantly lesser PCF matching error. In summary, the proposed improvements to sample synthesis enables unprecedented capabilities in exploratory sampling.

We demonstrate that using experiments in predictive modeling and hyper-parameter optimization.

V. EMPIRICAL STUDY: SAMPLING FOR PREDICTIVE MODELING

In this section, we study the qualitative performance of the proposed coverage-based design in predictive modeling, where one needs to recover unknown regression functions using a given set of sample observations. The goal of this study is to understand the impact of the improved coverage properties in the proposed design and the effectiveness of our sample synthesis algorithm. We consider both blind exploration, where the model is constructed only using one-shot exploratory samples, and sequential sampling, where the exploration samples are used to initialize a Bayesian optimization (Bayes-Opt) pipeline. Bayes-Opt [29] is a widely adopted sequential design framework typically employed for global optimization of complex functions. These methods begin by constructing a surrogate for the unknown function based on an initial sample, and then sequentially allocate the remaining design budget to quantify uncertainties of the surrogate, and utilize an acquisition function (e.g. expected improvement) to choose the next sample. We present comparisons to popular sampling methods, namely uniform random, Latin Hypercube sampling (LHS), Sobol (QMC) sequences, and SFSD, which is a state-of-the-art coverage-based design technique [19]. We show that the proposed approach produces superior recovery performance, thus establishing coverage-based designs as an effective solution for exploratory sampling.

Setup: We use the following benchmark functions from the global optimization literature: Alpine $N.1$ and Ackley in dimensions 3, 4 and 5, respectively. In order to evaluate the generalization of fitted functions, we generate 10^4 test samples using a regular grid in the sampling region, and use the mean squared error (MSE) with respect to the true function as the evaluation metric. For all experiments, we used random forest regressors with 100 trees, and the results reported were obtained by averaging over 20 independent realizations of sample designs.

Blind exploration: Figure 5 (a)-(f) compare the performance of our approach to the baseline methods in the fully exploratory case. It can be seen that the proposed design consistently outperforms popularly adopted sampling methods across varying N (50 to 200). Another striking observation is that there is significant variability in performance of the widely-adopted QMC sequences across dimensions, and as d increases it can sometimes perform even worse than uniform random samples. Furthermore, the poor performance of models learned using LHS and uniform random sampling for $d > 3$ can be directly attributed to their poor space-filling properties. Although SFSD and the proposed approach belong to the family of coverage-based designs, due to the improved coverage characteristics, our method consistently outperforms SFSD in all cases.

Sequential sampling: In this experiment, we study the impact of the choice for initial design on a Bayes-Opt pipeline. We consider an initial sampling budget of $N = 50$, and

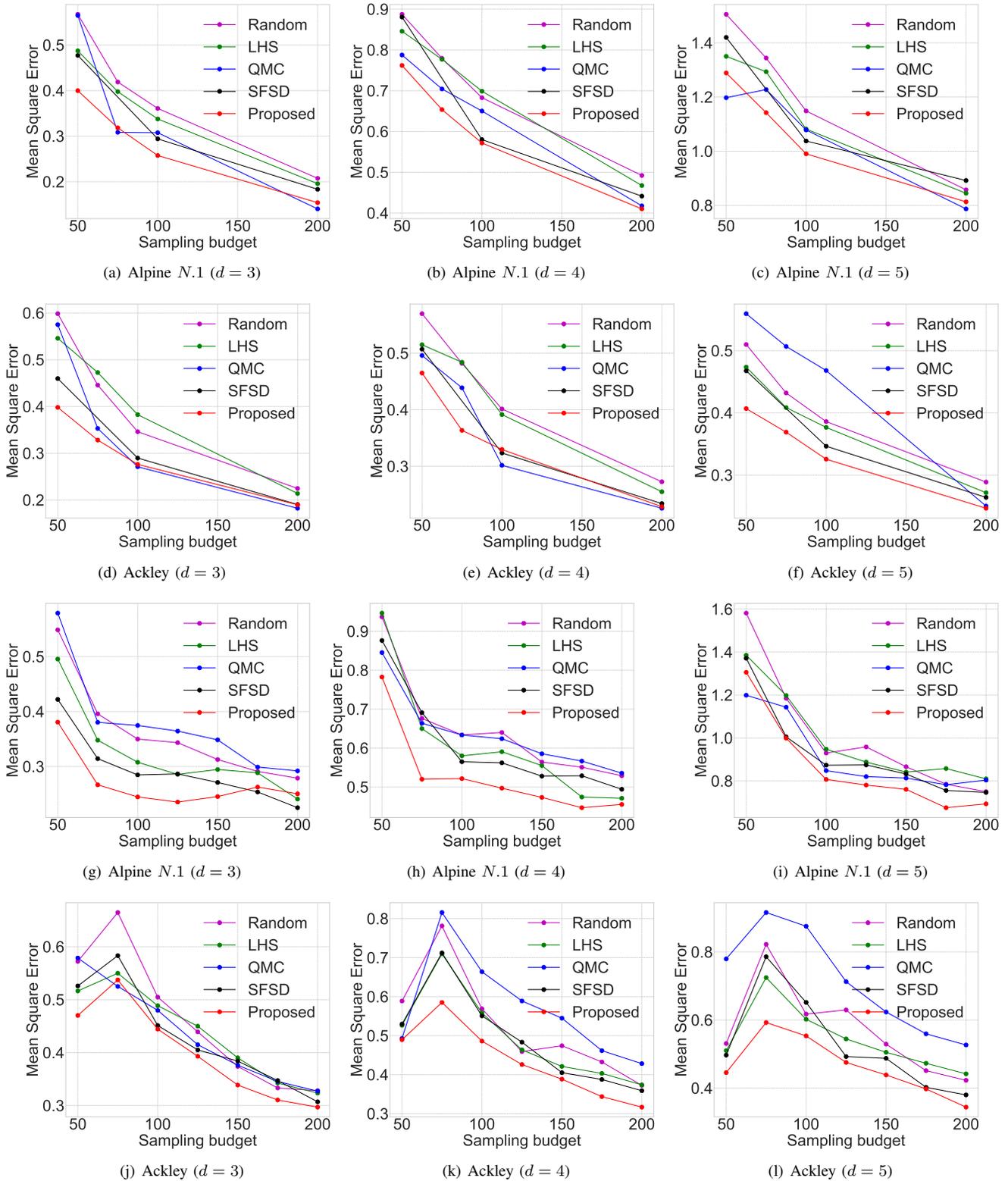


Figure 5: *Sampling for Predictive Modeling*: Performance of sample designs in recovering regression functions using : (a)-(f) blind exploration, (g)-(l) sequential sampling.

then sequentially sample 150 more samples to evaluate the behavior of Bayes-Opt on the same set of functions used in the previous case. Similar to the blind exploration case, we observe in Figure 5 (g)-(l) that the proposed design performs

significantly better compared to other state-of-the-practice choices. Although QMC sequences perform reasonably better than uniform random at $d = 3$, their performance degrades as d grows.

Table I: *Hyper-parameter search for a feature extractor-classifier pipeline: Average f1-score obtained over 10 realizations of exploratory sample designs for the 20-Newsdataset. This pipeline used tf-idf features and a SGD classifier.*

Blind Exploration						
d	N	QMC	LHS	Random	SFSD	Proposed
5	50	84.0995 ± 0	83.838 ± 0.415	83.872 ± 0.257	84.328 ± 0.166	84.441 ± 0.064
5	75	84.1758 ± 0	84.179 ± 0.237	84.196 ± 0.172	84.431 ± 0.138	84.414 ± 0.100
5	100	84.1543 ± 0	84.227 ± 0.189	84.161 ± 0.205	84.331 ± 0.127	84.463 ± 0.075
5	125	84.3501 ± 0	84.292 ± 0.131	84.286 ± 0.158	84.458 ± 0.073	84.482 ± 0.089
5	150	84.3294 ± 0	84.314 ± 0.141	84.336 ± 0.185	84.453 ± 0.058	84.527 ± 0.028

Table II: *Hyper-parameter search to build deep networks for MNIST digit classification: Best test accuracy obtained through the inclusion of hyper-parameter optimization using different sample designs. Note that, we consider both blind exploration and sequential sampling settings, and the results reported are averages over 10 independent realizations of the sample design.*

Blind Exploration						
d	N	QMC	LHS	Random	SFSD	Proposed
3	50	98.57 ± 0	91.198 ± 3.492	98.554 ± 0.229	98.691 ± 0.316	98.79 ± 0.198
3	100	98.82 ± 0	98.794 ± 0.171	98.688 ± 0.431	98.873 ± 0.124	98.896 ± 0.098
3	200	98.92 ± 0	98.932 ± 0.033	98.921 ± 0.126	98.975 ± 0.125	98.969 ± 0.035
4	50	98.66 ± 0	98.116 ± 0.690	97.818 ± 1.285	98.623 ± 0.325	98.806 ± 0.138
4	100	98.61 ± 0	98.70 ± 0.215	98.654 ± 0.216	98.748 ± 0.202	98.976 ± 0.157
4	200	98.199 ± 0	98.832 ± 0.117	98.902 ± 0.134	98.921 ± 0.075	98.932 ± 0.061
5	50	98.188 ± 0	97.996 ± 0.737	91.06 ± 3.485	98.622 ± 0.238	98.832 ± 0.134
5	100	98.77 ± 0	98.802 ± 0.163	98.642 ± 0.233	98.846 ± 0.188	98.834 ± 0.148
5	200	98.73 ± 0	98.992 ± 0.102	98.862 ± 0.131	98.944 ± 0.118	98.967 ± 0.068
Sequential Sampling						
d	N	QMC	LHS	Random	SFSD	Proposed
3	100	97.354 ± 0	97.466 ± 0.081	97.49 ± 0.445	97.367 ± 0.371	97.626 ± 0.128
4	100	97.581 ± 0	96.952 ± 0.562	97.492 ± 0.135	97.628 ± 0.198	97.597 ± 0.196
5	100	94.222 ± 0	97.296 ± 0.434	96.171 ± 0.950	97.487 ± 0.3134	97.662 ± 0.110

VI. APPLICATION: HYPER-PARAMETER SEARCH

Hyper-parameter search is critical to modern machine learning algorithms and resource-efficient optimization is directly linked to the scalability of the solutions. For this experiment, we consider both a conventional feature extractor-classifier pipeline and end-to-end deep learning systems, where the goal is to minimize the validation error [20]. The search space is characterized by a sparse set of locally optimal solutions, and requires effective sampling to rapidly choose a well-performing configuration. The evaluation metric that we use is the *precision*, i.e., the number of selected configurations that produces validation accuracies greater than a pre-defined threshold τ . We use this proxy metric [6], [20] since the global optimum is unknown, and more importantly identifying multiple locally optimal configurations in the search space reflects the ability of a sampling technique in characterizing the response surface. For completeness, we also include the widely-used *best validation accuracy* achieved over multiple realizations of the considered sample designs.

A. Conventional Feature Extractor-Classifer Pipeline

In this experiment, we consider the problem of choosing hyper-parameters for feature extraction and classification of

text documents in the *20-Newsdataset*. This is a collection of approximately 20,000 documents, partitioned evenly across 20 different newsgroups. We use a feature extractor-classifier pipeline that consists of : (a) *Count Vectorizer*: converts a collection of text documents into a matrix of token counts; (b) *Tf-idf-Transformer*: transforms a count matrix into a normalized (term-frequency times inverse document-frequency) tf-idf representation; and (c) linear classifier with stochastic gradient descent (SGD) training. We considered 5 hyper-parameters – document frequency threshold and maximum number of features in the feature extraction step, and 3 settings for classifier design : number of iterations, learning rate and regularization penalty.

We vary the sampling budget N in the range (50, 150) and compute the *f1-score* (macro-averaged) from the best performing configuration in the exploratory sample. We report the mean and standard deviation obtained using 10 independent realizations of the samples. As showed in Table I, in most of the cases, coverage-based designs outperform other random sampling baselines both in terms of expected performance and variance. More specifically, the proposed approach identifies the best configuration in every case. The superior performance of the proposed design over SFSD can be attributed to the

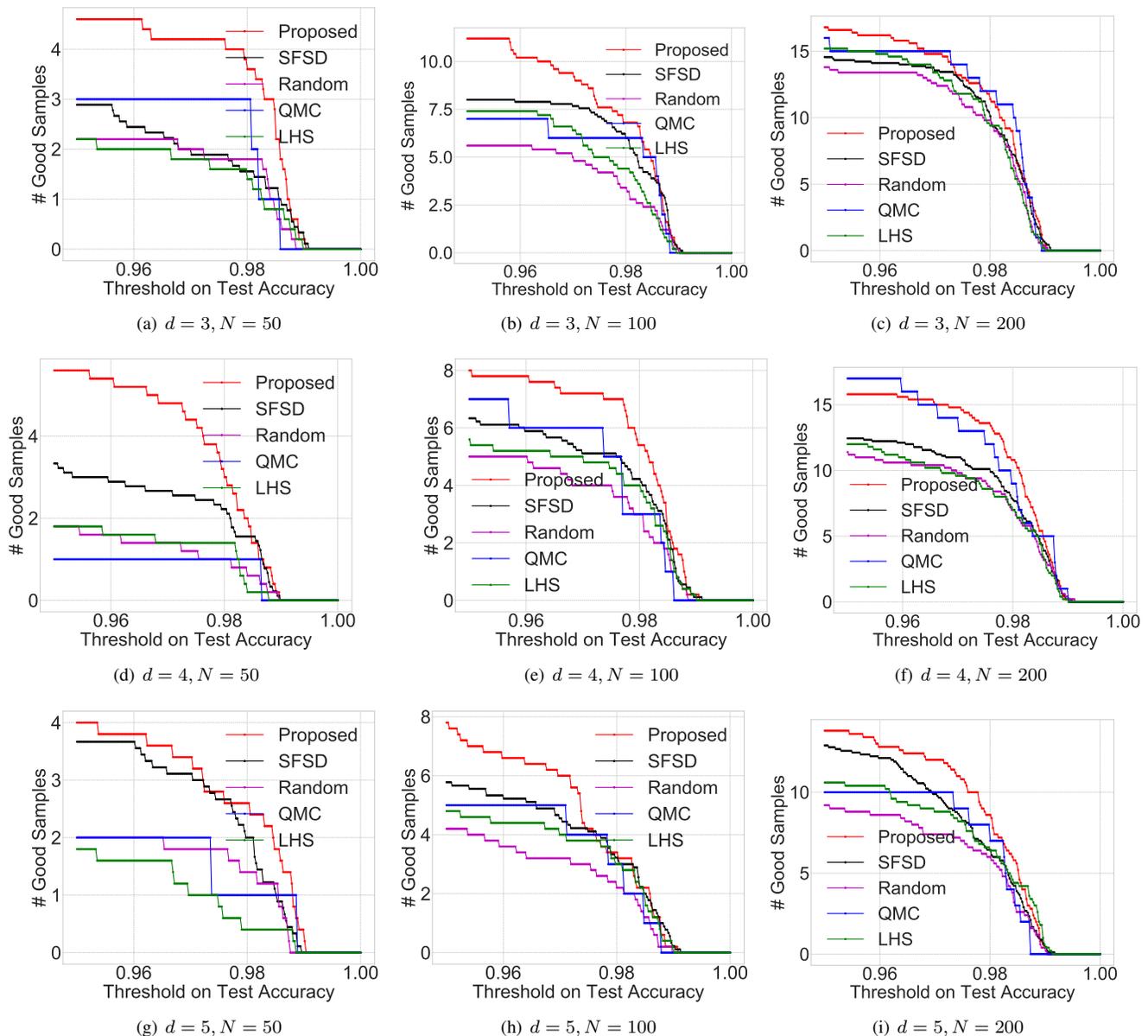


Figure 6: *Hyper-parameter search to build deep networks MNIST digit recognition: Precision metric obtained through blind exploration with different sample designs.*

improved coverage in the synthesized samples. On the other hand, conventional methods such as *LHS* and *Random* suffer from high variance across realizations.

B. Building Deep Models for MNIST Digit Classification

In this section, we consider the problem of building deep networks for classifying handwritten digits from MNIST, which contains 50,000 train and 10,000 test images. We evaluate the proposed sample design approach under blind exploration and sequential sampling settings.

Blind exploration: We use a simple CNN architecture: $conv[3 \times 3 \times 8] \rightarrow conv[3 \times 3 \times 16] \rightarrow FC[128] \rightarrow FC[64] \rightarrow FC[10]$. with ReLU activation and dropouts in after every layer. The training was carried out using gradient

descent with the momentum optimizer. The set of 5 hyper-parameters included learning rate, momentum and dropouts at the 2^{nd} *conv* layer, the 1^{st} *FC* layer and the 2^{nd} *FC* layer respectively. We also considered $4-d$ and $3-d$ subsets, where some of the dropouts fixed at 0.5. For this experiment, we used the sampling budgets $N = \{50, 100, 200\}$. In each case, we estimated the *precision* metric by varying the threshold τ , between 0.95 and 1. All results reported were averaged over 10 independent realizations.

Figure 6 and Table II illustrate the validation performance of different sample designs for this problem. We observe that the proposed design consistently achieves superior precision over existing experimental designs, thus ensuring a high probability of obtaining a generalizable model, particularly at lower sampling budgets. Although *LHS* and *QMC* samples perform

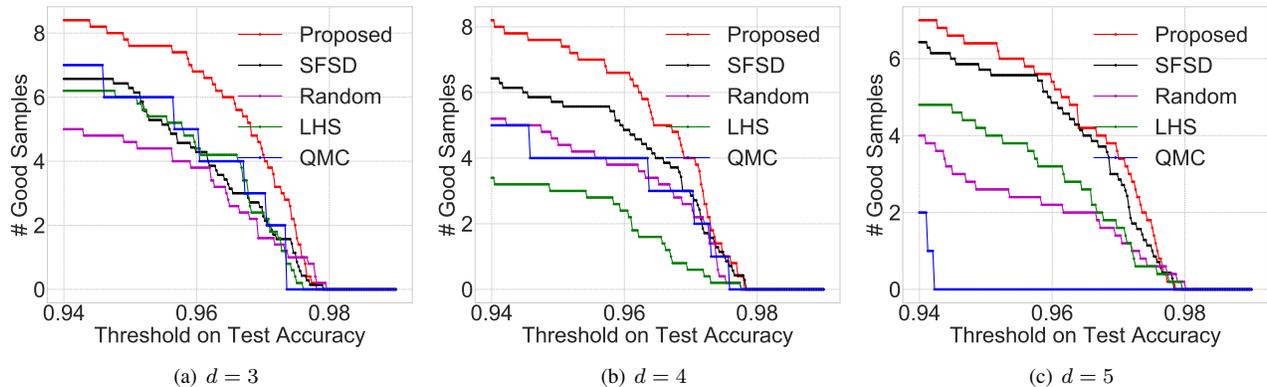


Figure 7: *Hyper-parameter search to build deep networks MNIST digit recognition*: Precision metric obtained through Bayes-Opt with different initial exploratory samples.

reasonably well in some cases, their performance degrades as d grows. Through the improved coverage characteristics, our approach sometimes identifies even twice as many local optima (based on the precision metric), thus motivating its use as an initializer for subsequent exploitation using Bayesian optimization. The results in Table II show that our method is able to sample the region of maximum interest consistently with low variance. Note that, the state-of-the-art SFSD design also performs consistently better than *Random* and discrepancy based designs in terms of test accuracy, but often demonstrates a larger variance.

Sequential Sampling: The success of Bayesian optimization relies on its ability to exploit uncertainties in the search space, to trade-off between exploration and exploitation. We argue that the choice of initial space-filling design can significantly impact the performance of sequential optimization in hyper-parameter search. For this experiment, we train a DNN architecture comprised of only dense layers: $FC[784] \rightarrow FC[512] \rightarrow FC[256] \rightarrow FC[64] \rightarrow FC[10]$, with ReLU activation and dropout after every layer. We used the same set of hyper-parameters as in the previous case, i.e., learning rate, momentum and dropout ratios. Experiments were conducted with an initial sampling budget of $N = 50$ samples and an additional 50 samples from sequential sampling, in dimensions 3, 4 and 5, respectively. Figure 7 and Table II demonstrate the impact of different initial exploratory samples on the sequential optimization performance. The gains over discrepancy-based and random designs is even more significant in this case, thus emphasizing coverage as a desired characteristic of exploratory designs. The consistency of our approach in its performance across dimensions, evidences its robustness when compared to other widely adopted designs.

C. Building CNN for Cifar-10 Image Classification

In this final experiment, we consider a $5 - d$ hyper-parameter search to train a CNN for classifying images from the CIFAR-10 dataset. The architecture used is as follows: $conv[3 \times 3 \times 32] \rightarrow conv[3 \times 3 \times 32] \rightarrow conv[3 \times 3 \times 64] \rightarrow conv[3 \times 3 \times 64] \rightarrow FC[512] \rightarrow FC[128] \rightarrow FC[10]$. with ReLU activation, max-pooling, and batch normalization after every convolution layer. Dropouts are included after the 2^{nd}

$conv$ layer, the 4^{th} $conv$ layer and the 1^{st} FC layer. The set of 5 hyper-parameters included learning rate, momentum and the 3 dropout ratios. For blind exploration, we used sampling budgets of $N = 50$ and $N = 100$. In case of sequential sampling, experiments were conducted with an initial budget of $N = 50$ samples and an additional 50 were sampled sequentially using Bayesian optimization. We report the mean and standard deviation of the best test accuracy achieved over 10 realizations in Table III and the precision metric in Figure 8. In all the cases, the proposed method achieves the best expected generalization performance. Along with the other experiments, this observation clearly strengthens the premise that coverage-based designs are highly effective for hyper-parameter search when compared to existing random sampling and discrepancy-based designs.

VII. CONCLUSIONS

We considered the problem of designing high quality exploratory samples. We introduced improved coverage Poisson disk sample designs using pair correlation function. We also proposed an approach to automatically determine the optimal parameters of the PDS designs. To generate these samples with high accuracy, we proposed an adaptive learning rate based gradient descent approach and showed that it significantly outperforms baseline methods. Finally, we evaluated the performance of PDS designs on predictive modeling and hyper-parameter search applications in both blind exploration and sequential search with Bayesian optimization. Experimental results show that the proposed PDS approach consistently outperforms state-of-the-art techniques, especially with low sampling budget.

VIII. ACKNOWLEDGMENTS

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process

Table III: *Hyper-parameter search to build CNNs for Cifar-10 image classification: Best test accuracy obtained through the inclusion of hyper-parameter optimization using different sample designs. Note that, we consider both blind exploration and sequential sampling settings, and the results are averages over 10 independent realizations of the sample design.*

Blind Exploration						
d	N	QMC	LHS	Random	SFSD	Proposed
5	50	80.36 ± 0	80.023 ± 0.393	80.217 ± 0.217	80.145 ± 0.409	80.522 ± 0.334
5	100	80.70 ± 0	80.338 ± 0.245	80.448 ± 0.236	80.488 ± 0.136	80.959 ± 0.374
Sequential Sampling						
d	N	QMC	LHS	Random	SFSD	Proposed
5	100	80.842 ± 0.0	80.436 ± 0.206	80.623 ± 0.561	80.866 ± 0.341	81.033 ± 0.294

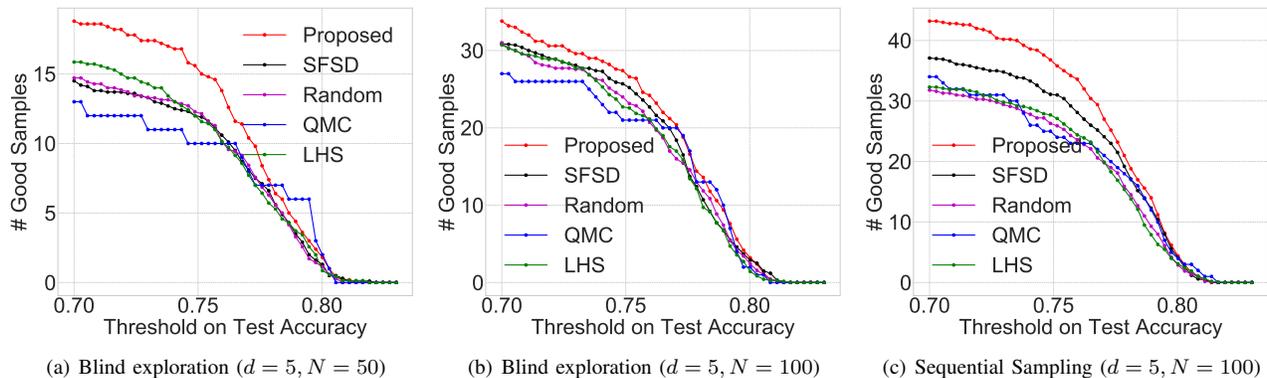


Figure 8: *Hyper-parameter search to build CNNs for Cifar-10 image classification: Precision metric obtained through blind exploration and Bayes-Opt with different initial exploratory samples.*

disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.

REFERENCES

- [1] R. A. Fisher, “The design of experiments. 1935,” *Oliver and Boyd, Edinburgh*, 1935.
- [2] H.-H. Bock and E. Diday, *Analysis of symbolic data: exploratory methods for extracting statistical information from complex data*. Springer Science & Business Media, 2012.
- [3] D. Csiba and P. Richtárik, “Importance sampling for minibatches,” *The Journal of Machine Learning Research*, vol. 19, pp. 962–982, 2018.
- [4] C.-Y. Wu, R. Manmatha, A. J. Smola, and P. Krahenbuhl, “Sampling matters in deep embedding learning,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2840–2848.
- [5] J. Bergstra and Y. Bengio, “Random search for hyper-parameter optimization,” *Journal of Machine Learning Research*, vol. 13, no. Feb, pp. 281–305, 2012.
- [6] O. Bousquet, S. Gelly, K. Kurach, O. Teytaud, and D. Vincent, “Critical hyper-parameters: No random, no cry,” *arXiv preprint arXiv:1706.03200*, 2017.
- [7] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [8] J. Asmuth, L. Li, M. L. Littman, A. Nouri, and D. Wingate, “A bayesian sampling approach to exploration in reinforcement learning,” in *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence*. AUAI Press, 2009, pp. 19–26.
- [9] S. Milli, P. Abbeel, and I. Mordatch, “Interpretable and pedagogical examples,” *arXiv preprint arXiv:1711.00694*, 2017.
- [10] A. K. Gupta, K. G. Smith, and C. E. Shalley, “The interplay between exploration and exploitation,” *Academy of management journal*, vol. 49, no. 4, pp. 693–706, 2006.
- [11] A. Doucet, S. Godsill, and C. Andrieu, “On sequential monte carlo sampling methods for bayesian filtering,” *Statistics and computing*, vol. 10, no. 3, pp. 197–208, 2000.
- [12] M. S. Ebeida, S. A. Mitchell, A. Patney, A. A. Davidson, and J. D. Owens, “A simple algorithm for maximal poisson-disk sampling in high dimensions,” *Computer Graphics Forum*, vol. 31, pp. 785–794, 2012.
- [13] M. S. Ebeida, A. Patney, S. A. Mitchell, K. R. Dalbey, A. A. Davidson, and J. D. Owens, “K-d darts: Sampling by k-dimensional flat searches,” *ACM Trans. Graph.*, vol. 33, no. 1, pp. 3:1–3:16, Feb. 2014.
- [14] A. B. Owen, “Monte carlo and quasi-monte carlo for statistics,” *Monte Carlo and Quasi-Monte Carlo Methods 2008*, pp. 3–18, 2009.
- [15] H. Niederreiter, *Random number generation and quasi-Monte Carlo methods*. Siam, 1992, vol. 63.
- [16] A. C. Öztireli and M. Gross, “Analysis and synthesis of point distributions based on pair correlation,” *ACM Trans. Graph.*, vol. 31, no. 6, pp. 170:1–170:10, Nov. 2012.
- [17] B. Kailkhura, J. J. Thiagarajan, P. T. Bremer, and P. K. Varshney, “Theoretical guarantees for poisson disk sampling using pair correlation function,” in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, March 2016, pp. 2589–2593.
- [18] B. Kailkhura, J. J. Thiagarajan, P.-T. Bremer, and P. K. Varshney, “Stair blue noise sampling,” *ACM Trans. Graph.*, vol. 35, no. 6, pp. 248:1–248:10, Nov. 2016. [Online]. Available: <http://doi.acm.org/10.1145/2980179.2982435>
- [19] B. Kailkhura, J. J. Thiagarajan, C. Rastogi, P. K. Varshney, and P.-T. Bremer, “A spectral approach for the design of experiments: Design, analysis and algorithms,” *Journal of Machine Learning Research (to appear)*, 2018.

- [20] J. S. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl, "Algorithms for hyper-parameter optimization," in *Advances in neural information processing systems*, 2011, pp. 2546–2554.
- [21] M. A. Z. Dippé and E. H. Wold, "Antialiasing through stochastic sampling," *SIGGRAPH Comput. Graph.*, vol. 19, no. 3, pp. 69–78, Jul. 1985. [Online]. Available: <http://doi.acm.org/10.1145/325165.325182>
- [22] R. L. Cook, "Stochastic sampling in computer graphics," *ACM Trans. Graph.*, vol. 5, no. 1, pp. 51–72, Jan. 1986. [Online]. Available: <http://doi.acm.org/10.1145/7529.8927>
- [23] J. Yellott, "Spectral consequences of photoreceptor sampling in the rhesus retina," *Science*, vol. 221, no. 4608, pp. 382–385, 1983. [Online]. Available: <http://science.sciencemag.org/content/221/4608/382>
- [24] M. S. Ebeida, A. A. Davidson, A. Patney, P. M. Knupp, S. A. Mitchell, and J. D. Owens, "Efficient maximal poisson-disk sampling," *ACM Trans. Graph.*, vol. 30, no. 4, pp. 49:1–49:12, Jul. 2011.
- [25] D. Heck, T. Schlömer, and O. Deussen, "Blue noise sampling with controlled aliasing," *ACM Trans. Graph.*, vol. 32, pp. 1–12, Jul. 2013.
- [26] M. N. Gamito and S. C. Maddock, "Accurate multidimensional poisson-disk sampling," *ACM Trans on Graphics (TOG)*, vol. 29, p. 8, 2009.
- [27] E. W. Weisstein, "Hypersphere packing." *From MathWorld – A Wolfram Web Resource*. [Online]. Available: <http://mathworld.wolfram.com/HyperspherePacking.html>
- [28] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [29] J. Snoek, H. Larochelle, and R. P. Adams, "Practical bayesian optimization of machine learning algorithms," in *Advances in neural information processing systems*, 2012, pp. 2951–2959.