

DETERMINANTAL POINT PROCESSES IMPLICITLY REGULARIZE SEMI-PARAMETRIC REGRESSION PROBLEMS*

MICHAËL FANUEL[†], JOACHIM SCHREURS[†], AND JOHAN A.K. SUYKENS[†]

Abstract. Semi-parametric regression models are used in several applications which require comprehensibility without sacrificing accuracy. Typical examples are spline interpolation in geophysics, or non-linear time series problems, where the system includes a linear and non-linear component. We discuss here the use of a finite Determinantal Point Process (DPP) for approximating semi-parametric models. Recently, Barthelmé, Tremblay, Usevich, and Amblard introduced a novel representation of some finite DPPs. These authors formulated *extended L-ensembles* that can conveniently represent partial-projection DPPs and suggest their use for optimal interpolation. With the help of this formalism, we derive a key identity illustrating the implicit regularization effect of determinantal sampling for semi-parametric regression and interpolation. Also, a novel *projected* Nyström approximation is defined and used to derive a bound on the expected risk for the corresponding approximation of semi-parametric regression. This work naturally extends similar results obtained for kernel ridge regression.

Key words. determinantal point processes, semi-parametric regression, Nyström approximation, implicit regularization

1. Introduction. Kernel methods provide a theoretically grounded framework for non-parametric regression and have been able to achieve excellent performance [51, 52] in the last years. In applications that require more explainability, a parametric component, usually a polynomial, is added to the kernel regressor. This semi-parametric model has the best of both worlds, a parametric component that is understandable for the user and a non-parametric kernel component that boosts the accuracy of the prediction. Full-size kernel regression problems do not scale well with the size of data sets, for that reason several approximations have been studied. In particular, in the case of massive data sets, smart sampling and sketching methods have allowed to scale up kernel ridge regression [43], while preserving its statistical guarantees. Not only to reduce memory requirements, sampling methods are interesting to reduce the number of parameters of such models for enhancing prediction speed, for instance, in the context of embedded applications. In this paper, we consider the specific setting of semi-parametric regression which generalizes and improves the interpretability of kernel ridge regression for the applications where a parametric (e.g., polynomial) estimator can be an educated guess. We combine this semi-parametric approach with a custom sampling scheme based on Determinantal Point Processes, thereby allowing to obtain subsets of important and diverse points. This leads to similar results to the ones obtained for kernel ridge regression [30], that is to say, DPP sampling implicitly regularizes (semi-parametric) kernel regression problems.

Sampling with a determinantal point process. Discrete Determinantal Point Processes (DPPs) provide elegant ways to sample random subsets $\mathcal{C} \subseteq \{1, \dots, n\}$, sometimes called ‘coresets’ [57], so that the selected items are diverse. In a word, discrete DPPs are represented by a marginal kernel, that is, a $n \times n$ matrix P with eigenvalues within $[0, 1]$, giving the inclusion probabilities: if \mathcal{C} is a random subset distributed according to a DPP with marginal kernel P , then the inclusion probabilities are

$$\Pr(\mathcal{E} \subseteq \mathcal{C}) = \det(P_{\mathcal{E}\mathcal{E}}),$$

where $P_{\mathcal{E}\mathcal{E}}$ is the square submatrix obtained by selecting the rows and columns of P indexed by \mathcal{E} . The off-diagonal entries of the marginal kernel are interpreted as similarity scores. Thus, a subset with a large probability is diverse. This can intuitively be seen thanks to the interpretation of the determinant of a positive definite matrix in terms of squared volume. In general, the expression of the probability for sampling a given subset $\Pr(\mathcal{C})$ is known but non trivial. Therefore, it is often very convenient to work with a L -ensemble,

*

Funding: EU: The research leading to these results has received funding from the European Research Council under the European Union’s Horizon 2020 research and innovation program / ERC Advanced Grant E-DUALITY (787960). This paper reflects only the authors’ views and the Union is not liable for any use that may be made of the contained information. Research Council KU Leuven: Optimization frameworks for deep kernel machines C14/18/068 Flemish Government: FWO: projects: GOA4917N (Deep Restricted Kernel Machines: Methods and Foundations), PhD/Postdoc grant This research received funding from the Flemish Government (AI Research Program). Ford KU Leuven Research Alliance Project KUL0076 (Stability analysis and performance improvement of deep reinforcement learning algorithms) EU H2020 ICT-48 Network TAILOR (Foundations of Trustworthy AI - Integrating Reasoning, Learning and Optimization). Leuven.AI Institute.

[†]KU Leuven, Department of Electrical Engineering (ESAT), STADIUS Center for Dynamical Systems, Signal Processing and Data Analytics, Kasteelpark Arenberg 10, B-3001 Leuven, Belgium. email: michael.fanuel@kuleuven.be

which is a DPP, denoted here by $DPP_L(L)$, such that

$$(1.1) \quad \Pr(\mathcal{C}) = \det(L_{\mathcal{C}\mathcal{C}}) / \det(\mathbb{I} + L),$$

where L is a $n \times n$ positive semi-definite matrix. The marginal kernel of an L -ensemble has the following simple expression

$$P = L(L + \mathbb{I})^{-1},$$

which is a matrix encountered in kernel ridge regression as we explain hereafter. In the context of *sketched* kernel ridge regression, sampling with L -ensemble DPPs yields very simple theoretical guarantees displaying an implicit regularization effect [30]. Before discussing semi-parametric regression, we briefly outline the known results about kernel ridge regression.

Sketching Kernel Ridge Regression. Given input-output pairs $(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$ for $1 \leq i \leq n$, kernel ridge regression (KRR) estimates a function of the form $f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i)$ with the help of a positive semi-definite kernel function $k(\mathbf{x}, \mathbf{x}')$ defined on $\mathbb{R}^d \times \mathbb{R}^d$, such as the Gaussian kernel $k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|_2^2 / \sigma^2)$. Classically, the numerical solution of KRR relies on a $n \times n$ positive semi-definite kernel matrix

$$K = [k(\mathbf{x}_i, \mathbf{x}_j)]_{1 \leq i, j \leq n},$$

constructed from the input data. Such a matrix can be potentially large if the size of the data set is large. Therefore, low rank approximations of K have been developed [59, 52, 30, 29]. In particular, an L -ensemble DPP can be used to select subsets $\mathcal{C} = \{c_1, \dots, c_k\}$ of $\{1, \dots, n\}$ in order to sample a subset of entries of K . In this context, a natural choice is $L = K/\lambda$ for some $\lambda > 0$. Then, it is customary to approximate K thanks to the low rank Nyström method which uses its submatrices, such as the square $k \times k$ submatrix $K_{\mathcal{C}\mathcal{C}}$. Sampling is conveniently done with the use of a $n \times k$ sampling matrix, that is obtained by selecting the columns of the identity matrix indexed by \mathcal{C} as follows $C = (\mathbf{e}_{c_1} \cdots \mathbf{e}_{c_k})$, where \mathbf{e}_i denotes the i -th element of the canonical basis. In the case of the Nyström approximation, one considers the pseudo-inverse of the *sparse* matrix $CK_{\mathcal{C}\mathcal{C}}C^\top$ (see below), which is a $n \times n$ matrix whose entry (i, j) is K_{ij} if $i, j \in \mathcal{C}$ and zero otherwise. Explicitly, the *common* Nyström approximation of K is defined as

$$(Nyström) \quad K(CK_{\mathcal{C}\mathcal{C}}C^\top)^+K = KCK_{\mathcal{C}\mathcal{C}}^+C^\top K,$$

where $(\cdot)^+$ denotes the Moore-Penrose pseudo-inverse¹. This subsampling with L -ensembles has an implicit regularization effect [30], which is based on the following expectation formula, also independently shown by [47]:

$$(1.2) \quad \mathbb{E}_{\mathcal{C}}(CK_{\mathcal{C}\mathcal{C}}C^\top)^+ = (K + \lambda\mathbb{I})^{-1},$$

where $\mathcal{C} \sim DPP_L(L)$ with $L = K/\lambda$ and $\lambda > 0$. Varying λ allows to vary the expected size of the subset and the amount of regularization. A similar identity has been revisited in the context of fixed-size L -ensemble DPP in [54]. Albeit the exact sampling of a L -ensemble DPP has a time complexity $\mathcal{O}(n^3)$, the obtained expected error for Nyström approximation $\mathbb{E}[K - KCK_{\mathcal{C}\mathcal{C}}^+C^\top K] = \lambda K(K + \lambda\mathbb{I})^{-1}$ for $\mathcal{C} \sim DPP_L(L)$ provides a generalization of error bounds obtained with Ridge Leverage Score sampling [26, 46, 51]. To the best of our knowledge, such results have not been obtained for semiparametric regression problems generalizing KRR.

Partial projection DPPs and semi-parametric regression. There are DPPs which are not L -ensembles, for instance, the projection DPPs for which the marginal kernel is a projector, i.e., a symmetric matrix such that $P^2 = P$. In practice, it is often convenient to have a simple formula for the probability that a subset is sampled, i.e., $\Pr(\mathcal{C})$. Therefore, the elegant framework of ‘extended L -ensembles’ has been introduced in [4] which provides a handy formula generalizing (1.1). This formalism is used extensively in this paper to deal with partial-projection DPPs. In Layman’s Terms, both partial projection DPPs and semi-parametric regression rely on mathematical expressions involving a sum of two objects living in orthogonal subspaces. This analogy is the main motivation to consider approximations of semi-parametric regression models with partial projection DPPs. For a partial-projection DPP, denoted by $DPP(L, V)$, the marginal kernel is of the following form

$$(1.3) \quad P = \mathbb{P}_V + \tilde{L}(\tilde{L} + \mathbb{I})^{-1},$$

¹See Lemma 6.7 in Appendix for a formal statement concerning the pseudo-inverse of this kind of matrices.

where the matrix V is a $n \times p$ matrix with full column rank and $\mathbb{P}_V = V(V^\top V)^{-1}V^\top$ is the projection on its column space, while the matrix $\tilde{L} = \tilde{K}/\lambda$ with $\lambda > 0$ is defined thanks to the $n \times n$ projected kernel

$$\tilde{K} \triangleq \mathbb{P}_{V^\perp} K \mathbb{P}_{V^\perp} \quad \text{with } \mathbb{P}_{V^\perp} = \mathbb{I} - \mathbb{P}_V.$$

In what follows, such a projected quantity is denoted by using a tilde. It is natural to assume \tilde{L} to be positive semi-definite so that the inverse matrix in (1.3) is well-defined.

Extended L -ensembles, that we describe below, represent partial-projection DPPs (see (1.3)) by giving a convenient formula for $\Pr(\mathcal{C})$. The reference [4] also points out a connection between extended L -ensembles and optimal interpolation in Section 2.8.2. This remark has motivated the following case study: the Nyström approximation of semi-parametric regression problem. The problem consists in recovering a function from noisy function values

$$y_i = z_i + \epsilon_i \quad \text{with } z_i = f(\mathbf{x}_i), \text{ and } 1 \leq i \leq n,$$

where ϵ_i denotes i.i.d. $\mathcal{N}(0, \sigma^2)$ noise. We consider here the semi-parametric model (see Figure 1 for an illustration)

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i) + \sum_{m=1}^p \beta_m p_m(\mathbf{x}),$$

where the first term is the non-parametric component associated to a *conditionally* positive semi-definite kernel $k(\mathbf{x}, \mathbf{x}')$, while the second term is the parametric component that is typically given by polynomials. The estimation problem amounts to solve $\hat{f} = \arg \min_f \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 + \gamma J(f)$, for a suitable penalization functional $J(f)$ defined in Section 3.3, and given $\gamma > 0$. Then, the marginal kernel (1.3) appears interestingly in the known formula for the in-sample estimate of a semi-parametric γ -regularized least squares problem,

$$\hat{z} = P \mathbf{y} \quad \text{with } P = \mathbb{P}_V + \tilde{L}(\tilde{L} + \mathbb{I})^{-1},$$

with $\tilde{L} = \tilde{K}/(n\gamma)$ and where \hat{z}_i denotes the estimated function value $z_i = f(\mathbf{x}_i)$ for $1 \leq i \leq n$ and $\gamma > 0$ is a regularization parameter. In this setting, the V and K matrices used in (1.3) are identified with the following matrices obtained from the parametric and non-parametric components:

$$V = [p_m(\mathbf{x}_i)]_{1 \leq i \leq n, 1 \leq m \leq p} \quad \text{and } K = [k(\mathbf{x}_i, \mathbf{x}_j)]_{1 \leq i, j \leq n}.$$

We now outline the contributions of this paper.

1.1. Contributions.

Sketched semi-parametric regression. First and foremost, a *key* contribution of this paper is a formula analogous to (1.2) involving a sampling with a custom partial-projection DPP. As it is explained above, full-fledged semi-parametric regression involves two orthogonal components, associated to the matrices V and $\mathbb{P}_{V^\perp} K \mathbb{P}_{V^\perp}$ respectively. To preserve this orthogonal decomposition for the sketched problem, we begin by defining a *sampling of the rows* of the matrix V —which stores the non-parametric component of the regression problem—and a sampling of rows and columns of K as follows

$$V_{\mathcal{C}} = [p_m(\mathbf{x}_i)]_{i \in \mathcal{C}, 1 \leq m \leq p} \quad \text{and } K_{\mathcal{C}\mathcal{C}} = [K_{ij}]_{i, j \in \mathcal{C}}.$$

We analyse a sketched regression problem which is constructed so that the orthogonality between the parametric and non-parametric components is preserved. A key ingredient to analyse this problem is the projector $\mathbb{P}_{V_{\mathcal{C}}^\perp}$ onto the orthogonal of the column space of $V_{\mathcal{C}} = C^\top V$. Then, we address the following question:

What is the relationship between the sketched regression problem associated to $V_{\mathcal{C}}$ and $\mathbb{P}_{V_{\mathcal{C}}^\perp} K_{\mathcal{C}\mathcal{C}} \mathbb{P}_{V_{\mathcal{C}}^\perp}$, and the full regression problem associated to V and $\mathbb{P}_{V^\perp} K \mathbb{P}_{V^\perp}$?

Our main result, Theorem 4.1 given hereafter, implies the following identity for the expectation of the pseudo-inverse²

$$(1.4) \quad \mathbb{E}_{\mathcal{C}} \left(C \mathbb{P}_{V_{\mathcal{C}}^\perp} K_{\mathcal{C}\mathcal{C}} \mathbb{P}_{V_{\mathcal{C}}^\perp} C^\top \right)^+ = \mathbb{P}_{V^\perp} (\tilde{K} + \lambda \mathbb{I})^{-1} \mathbb{P}_{V^\perp} \quad \text{with } \tilde{K} = \mathbb{P}_{V^\perp} K \mathbb{P}_{V^\perp},$$

²Technically, we require here that K is conditionally positive semi-definite with respect to V , i.e., $\mathbb{P}_{V^\perp} K \mathbb{P}_{V^\perp}$ is positive semi-definite.

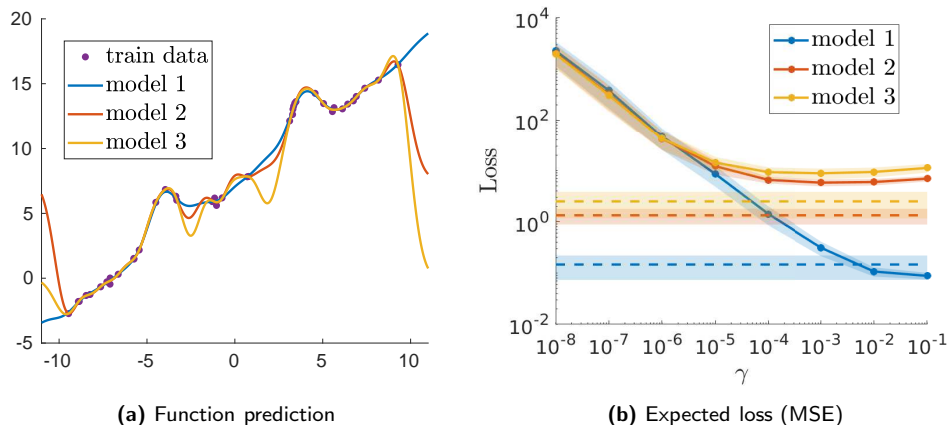


Fig. 1: A Toy example of semi-parametric regression with a Gaussian kernel ($d = 1$). Figure 1a shows the training points and the estimated function with $\sigma = 1$ and best performing γ . Only the semi-parametric model, i.e. model 1: $\hat{f}(x) = \beta_1 + \beta_2 x + \sum_i \alpha_i k(x, x_i)$, predicts the linear trend in low density regions as well outside the training interval $[-10, 10]$, contrary to model 2 (LS-SVM, i.e. $\beta_2 = 0$) and model 3 (KRR, i.e. $\beta_1 = \beta_2 = 0$). The MSE of each model with bandwidth $\sigma = 1$ is visualized as a function of the regularization parameter in Figure 1b. The dashed line shows the best performance when cross-validating over both γ and σ . See Section 3.2 for more details.

where \mathcal{C} is sampled according to the partial-projection $DPP(K/\lambda, V)$. We emphasize that there is no trivial connection between \mathbb{P}_{V^\perp} and $\mathbb{P}_{V_{\mathcal{C}}^\perp}$, which are the projectors onto the orthogonal of the column spaces of V and $V_{\mathcal{C}}$ respectively. The implicit regularization in (1.4) is ‘conditional’ since it occurs only within the subspace orthogonal to V . As in the case of L -ensembles, the real number $\lambda > 0$ influences the expected subset size and the amount of regularization. The identity (1.4) is novel to the best of our knowledge and is also instrumental to derive two key contributions of this paper.

Projected Nyström approximation. In Section 5, we define a projected Nyström approximation $\widetilde{L}(\mathcal{C})$ of the projected kernel matrix \widetilde{K} under the assumption that \widetilde{K} is positive semi-definite:

$$\text{(Projected Nyström)} \quad \widetilde{L}(\mathcal{C}) \triangleq \widetilde{K} S(\mathcal{C}) \left(S(\mathcal{C})^\top \widetilde{K} S(\mathcal{C}) \right)^+ S(\mathcal{C})^\top \widetilde{K},$$

where the sketching matrix is $S(\mathcal{C}) = CB(\mathcal{C}) \in \mathbb{R}^{n \times (k-p)}$, with $B(\mathcal{C}) \in \mathbb{R}^{k \times (k-p)}$ a matrix whose columns are an orthonormal basis of the orthogonal of the column space of $V_{\mathcal{C}}$. The projected Nyström naturally extends the common Nyström approximation to semi-parametric regression problems and is essential for scaling the model to larger data sets. The low rank approximation of the projected kernel matrix can be constructed conveniently with submatrices of the original kernel K . Indeed, it is not necessary to construct explicitly the sketching matrix $S(\mathcal{C})$. In comparison with the common Nyström approximation, the sketching matrix involves here a projection since $B(\mathcal{C})B(\mathcal{C})^\top = \mathbb{P}_{V_{\mathcal{C}}^\perp}$. Importantly, we give an expected error formula for the projected Nyström approximation in Corollary 5.6,

$$\mathbb{E}_{\mathcal{C}}[\widetilde{K} - \widetilde{L}(\mathcal{C})] = \lambda \widetilde{K} (\widetilde{K} + \lambda \mathbb{I})^{-1}, \text{ where } \mathcal{C} \sim DPP(K/\lambda, V).$$

Notice that the expected subset size of $\mathcal{C} \sim DPP(K/\lambda, V)$ is given by

$$\mathbb{E}_{\mathcal{C}}[|\mathcal{C}|] = p + d_{\text{eff}}(\widetilde{K}/\lambda), \text{ with } d_{\text{eff}}(\widetilde{K}/\lambda) = \text{Tr} \left(\widetilde{K} (\widetilde{K} + \lambda \mathbb{I}_n)^{-1} \right),$$

where p is the number of columns of V and with $\widetilde{K} = \mathbb{P}_{V^\perp} K \mathbb{P}_{V^\perp}$. The interpretation of the above identities is that a small $\lambda > 0$ yields a large number of samples and a small error on expectation. This extends similar results obtained independently in [30, Corollary 2] and [14].

Stability result. In Section 5.5, we give an expected risk bound for the estimator \hat{z}_N of the γ -regularized semi-parametric regression obtained with the projected Nyström approximation, that is,

$$\mathbb{E}_{\mathcal{C}} \left[\sqrt{\frac{\mathcal{R}(\hat{z}_N)}{\mathcal{R}(\hat{z})}} \right] \leq 1 + \frac{\lambda}{n\gamma} d_{\text{eff}}(\tilde{K}/\lambda), \text{ with } \mathcal{C} \sim \text{DPP}(K/\lambda, V),$$

where the expected risk of the estimator \hat{z} is $\mathcal{R}(\hat{z}) \triangleq \mathbb{E}_{\epsilon} \|\hat{z} - z\|_2^2$, as it is detailed in Theorem 5.8 hereafter. This stability result indicates that the estimation thanks to the Nyström approximation cannot be arbitrarily worse than the estimation obtained without approximation.

Two different applications are considered within the penalized kernel regression framework. 1) The first case occurs when the output values (the y_i 's) are initially unknown to the user and costly to retrieve. This is for example the case in an active learning approach where the data points have to be manually labelled or when measurements are expensive. This application is known as 'discrete' experimental design and was previously studied, e.g., for linear regression in [19, 17, 13]. In this setting, one interpolates on a small number of selected landmark points to minimize the number of necessary labeled points. The question now poses itself: what is a good way of selecting points such that the performance is maintained together with a good conditioning of the linear system? In this paper, we propose a determinantal design approach. 2) The user has knowledge of the full response vector \mathbf{y} , but the (embedded) application requires a number of parameters smaller than $n + p$, or the number of data points n is too large to solve the corresponding linear system.

Random design regression. Incidentally, we provide in Section 5.2 a discrete random design method for parametric problems of the type $\min_{\beta} \|V\beta - \mathbf{y}\|_2^2$. Essentially, a partial projection DPP is used in order to sample a subset \mathcal{C} so that the estimator $\hat{\beta}(\mathcal{C}) = \arg \min_{\beta} \|V_{\mathcal{C}}\beta - \mathbf{y}_{\mathcal{C}}\|_2^2$ is unbiased. These result are analogous to those of [21] although the sampling algorithm is different. Our analysis directly follows from the main result in Theorem 4.1, and provides an alternative method generalizing volume sampling which might be of independent interest.

1.2. Related work. While it is currently an active topic of research in the context of deep neural networks, implicit regularization has been studied already previously in [39, 40]. Recently, DPPs and implicit regularization also appeared in the context of double descent phenomena [16], while we refer to [18] for a review. DPPs are useful methods to sample diverse subsets that have been applied in machine learning in variety of tasks, such as diverse recommendations, summarizing text or search tasks [37]. Implicit regularization is not specific to sampling, since it is also observed with Gaussian and Rademacher sketches of Gram matrices in [15], although a closed-form formula can be advantageously derived with L -ensemble sampling.

As it was mentioned earlier in the introduction, large scale KRR has been successfully solved thanks to the Nyström method combined with smart approximations of ridge leverage score (RLS) sampling in [43, 51] for data sets of several million points. Sampling with RLSs can be interpreted as an approximation of L -ensemble DPP sampling, where negative dependence is neglected [18, 30, 54]. In this spirit, new bounds on the Nyström approximation errors with L -ensemble sampling, naturally generalize the bounds obtained with RLS sampling [26]. The results of the aforementioned papers have been extended to Column Subset Selection Problems (CSSP) in [14].

Semi-parametric models are useful tools when some domain knowledge exists about the function to be estimated (e.g., a user wants to correct the data for a linear trend) or more understandability is required from a model [53, 36]. These models combine a parametric part which is easy to understand and non-parametric term to improve performance. Semi-parametric models are used in some critical applications where a user wants to have an understandable model, without sacrificing accuracy [27, 28, 55].

A natural application of conditionally positive semi-definite matrices is radial basis function interpolation [45], which is an attractive method for interpolating and smoothing function values on scattered points in the plane. However, a potential problem is that their computation involves the solution of a linear system that is often ill-conditioned for large data sets. Importantly, the paper [5] studies a slightly different question, which mainly concerns the case of thin-plate splines. Thin-plate spline basis functions have been shown to be very accurate in medical imaging [11], surface reconstruction [10], as well as other engineering applications [8]. Given a fixed set of landmark points, the authors of [5] propose an elegant method for choosing a suitable basis of the function space so that the linear system under study has an improved condition number. This strategy is also described in [58].

Random design regression has been studied recently from the viewpoint of repulsive point processes, whereas optimal design has received already a lot of attention in statistics (see for example [33] or [50] for a review). Prominent recent works using random designs involve volume sampling [21, 22] or use the Bayesian perspective [17]. Several extensions of DPP and volume sampling have been developed recently such as in [16] or in the generalization to Polish spaces in [49].

1.3. Organization of the paper. In Section 2, we introduce basic definitions. Then, the penalized semi-parametric regression and interpolation problems that we study in this paper are discussed in Section 3. There, we consider two case studies: thin-plate splines regression and Gaussian kernel semi-parametric regression. Next, in Section 4, we explain the implicit regularization effect of a determinantal design for optimal interpolation. In Section 5, a large scale semi-parametric regression problem and its approximation thanks to a determinantal sampling are discussed together with a custom Nyström approximation. We also provide the stability result for the expected risk of this approximation, which was announced in the introduction. Technical proofs and results are deferred to Appendix.

1.4. Notations. Matrices (A, B, K, \dots) are denoted by upper-case letters, whereas (column) vectors are denoted by bold lower-case letters, e.g., $\mathbf{x} = [x_1, \dots, x_d]^\top \in \mathbb{R}^{d \times 1}$. As mentioned above, the canonical basis of \mathbb{R}^n is written \mathbf{e}_i for $1 \leq i \leq n$. The constant $n \times 1$ vector of ones is $\mathbf{1}_n$. We write $A \succeq 0$ (resp. $A \preceq 0$) if A (resp. $-A$) is positive semi-definite (*psd*), while $A \succeq B$ indicates that $A - B$ is *psd*. The $n \times n$ identity matrix is written \mathbb{I}_n . The Moore-Penrose pseudo-inverse of a matrix A is denoted here by A^+ , and is given by $A^+ = (A^\top A)^{-1} A^\top$ if A has full column rank. We use calligraphic letters to denote sets, with the following exception $[n] = \{1, \dots, n\}$. In this paper, we consider the problem of sampling subsets $\mathcal{C} \subseteq [n]$. Then, it is convenient to use a sampling matrix $C \in \mathbb{R}^{n \times k}$ obtained by selecting the columns of the identity matrix corresponding to \mathcal{C} . The sampled submatrices are denoted as follows: $V_{\mathcal{C}} = C^\top V$ where $V \in \mathbb{R}^{n \times p}$ and $A_{\mathcal{C}\mathcal{C}} = C^\top A C$ with $A \in \mathbb{R}^{n \times n}$. The characteristic function of a set X is written $\mathbb{1}(X)$.

2. Extended L -ensembles and definitions . To begin, we recall some definitions that were introduced in [4]. First, an essential element is the non-negative pair, which is used to define the parametric and non-parametric components of the regressors.

DEFINITION 2.1 (non-negative pair [4]). *Let $V \in \mathbb{R}^{n \times p}$ be a matrix with full column rank, and let $A \in \mathbb{R}^{n \times n}$ be a conditionally positive semi-definite matrix with respect to V , i.e., a matrix satisfying $\tilde{A} \succeq 0$ where*

$$\tilde{A} \triangleq \mathbb{P}_{V^\perp} A \mathbb{P}_{V^\perp},$$

with \mathbb{P}_{V^\perp} the linear projector on the orthogonal of the space spanned by the columns of V . Then, the couple (A, V) is called a Non-Negative Pair (NNP).

The NNPs are closely related to a certain class of kernel functions, called conditionally positive semi-definite, whose kernel matrices are *psd* only within the orthogonal of a subspace. Several examples of these functions have been used in the context of optimal interpolation of scattered data.

DEFINITION 2.2 (conditionally positive semi-definite kernels). *A kernel function $k(\mathbf{x}, \mathbf{x}')$ is conditionally positive semi-definite with respect to $\{p_m(\mathbf{x})\}_{m=1}^p$, if for all finite sets $\{\mathbf{x}_i\}_{i=1}^n$ so that $V_{ij} = [p_j(\mathbf{x}_i)]$ is full column rank, the matrix $K_{ij} = [k(\mathbf{x}_i, \mathbf{x}_j)]_{i,j}$ is such that $\mathbb{P}_{V^\perp} K \mathbb{P}_{V^\perp} \succeq 0$.*

The following example of NNP is classical. Other examples will be considered in the context of thin-plate spline interpolation.

EXAMPLE 1. *The kernel $k(\mathbf{x}, \mathbf{x}') = -\|\mathbf{x} - \mathbf{x}'\|_2^2$ is a conditionally positive semi-definite kernel with respect to the constant function. Indeed, let $\mathbf{x}_i \in \mathbb{R}^d$ with $1 \leq i \leq n$. Then, the Gram matrix $K = [k(\mathbf{x}_i, \mathbf{x}_j)]_{i,j}$ satisfies $\mathbf{v}^\top K \mathbf{v} \geq 0$ for all $\mathbf{v} \in \mathbb{R}^d$ such that $\mathbf{1}_n^\top \mathbf{v} = 0$. Then, $(K, \mathbf{1}_n)$ is a NNP. Other examples are given by the generalized multiquadrics mentioned in [4].*

Equipped with these definitions, a DPP associated to a NNP can be formally defined thanks to the following representation called extended L -ensemble.

DEFINITION 2.3 (extended L -ensemble [4]). *Let (L, V) be a NNP and $\tilde{L} = \mathbb{P}_{V^\perp} L \mathbb{P}_{V^\perp}$. Then, an*

extended L -ensemble $\mathcal{Y} \sim DPP(L, V)$ satisfies

$$\Pr(\mathcal{Y} = \mathcal{C}) = N^{-1} \times \det \begin{pmatrix} L_{\mathcal{C}\mathcal{C}} & V_{\mathcal{C}} \\ V_{\mathcal{C}}^{\top} & 0 \end{pmatrix},$$

where the normalization writes $N = (-1)^p \det(\tilde{L} + \mathbb{I}) \det(V^{\top} V)$.

The connection between the marginal kernel of a partial projection DPP given in (1.3) and extended L -ensembles is discussed in [4]. It is worth emphasizing that extended L -ensemble elegantly combine the probability mass functions of projection DPPs and L -ensemble (cfr. (1.1)). Indeed, the determinant of the block matrix in Definition 2.3 can be conveniently calculated (see Lemma 6.10 in Appendix), so that an equivalent expression writes

$$(2.1) \quad \Pr(\mathcal{Y} = \mathcal{C}) = \frac{\det \left(B(\mathcal{C})^{\top} \tilde{L}_{\mathcal{C}\mathcal{C}} B(\mathcal{C}) \right) \det(V_{\mathcal{C}}^{\top} V_{\mathcal{C}})}{\det(\mathbb{I} + \tilde{L}) \det(V^{\top} V)},$$

where, as defined in the introduction, $B(\mathcal{C}) \in \mathbb{R}^{k \times (k-p)}$ is a matrix whose columns form an orthonormal basis of the orthogonal of the column space of $V_{\mathcal{C}}$ and conveniently can be calculated using the QR decomposition.

Remark 2.4 (limit case). Consider the partial projection process $DPP(t\mathbb{I}, V)$ with $t > 0$. Its probability mass function converges pointwisely to

$$\Pr(\mathcal{Y} = \mathcal{C}) = \mathbb{1}(|\mathcal{C}| = p) \frac{\det(V_{\mathcal{C}}^{\top} V_{\mathcal{C}})}{\det(V^{\top} V)},$$

as $t \rightarrow 0$. Thus, the limit process is actually a projection DPP, i.e., a DPP with the following marginal kernel $P = \mathbb{P}_V \triangleq V(V^{\top} V)^{-1} V^{\top}$.

The subset size of an extended L -ensemble is also a random variable. Explicitly, the expected size of $\mathcal{C} \sim DPP(K/\lambda, V)$ is

$$\mathbb{E}[|\mathcal{C}|] = \text{Tr} \left(\tilde{K}(\tilde{K} + \lambda\mathbb{I})^{-1} \right) + p.$$

The parameter $\lambda > 0$ allows to vary the sample size, i.e., a small λ value yields a large sample size on expectation and conversely. This can be shown thanks to the marginal kernel of an extended L -ensemble, which was given in (1.3).

Remark 2.5 (sampling). The sampling algorithm used in this paper relies on Algorithm 3 in [57] (see also [37]). In all the numerical simulations, we use a fixed-size DPP sampling. A fixed-size DPP is a DPP conditioned on a fixed subset size. Our choice is motivated by the asymptotic equivalence between DPPs and fixed-size DPPs [3]. Importantly, the number of operations to exactly sample a DPP is $\mathcal{O}(n^3)$. This cost can be reduced if the marginal kernel has a low rank structure. Nonetheless, several approximation techniques have been published recently [9, 12] in order to alleviate the cost of sampling fixed-size DPPs, especially if the subset size is small.

Remark 2.6 (leverage scores). Leverage scores [24] and λ -ridge leverage scores [26] have been designed for randomized matrix approximations with i.i.d. sampling. Ridge leverage score (RLS) sampling can be seen as an approximation of DDP sampling by neglecting the negative dependence [18, 30]. In regards to the marginal kernel (1.3), the marginal probabilities of a partial-projection $\mathcal{Y} \sim DPP(\tilde{K}/\lambda, V)$ are the sum of a leverage score and a λ -ridge leverage score:

$$\Pr(i \in \mathcal{Y}) = \underbrace{\mathbf{e}_i^{\top} V(V^{\top} V)^{-1} V^{\top} \mathbf{e}_i}_{\text{leverage score}} + \underbrace{\mathbf{e}_i^{\top} \tilde{K}(\tilde{K} + \lambda\mathbb{I})^{-1} \mathbf{e}_i}_{\text{ridge leverage score}}.$$

Hence, partial projection DPP sampling provide a generalization of leverage score sampling. To the best of our knowledge, the above combination of leverage scores and λ -ridge leverage score sampling has received up to now little interest in the literature.

3. Basics of penalized semi-parametric regression. We begin by introducing the framework of semi-parametric regression with a *psd* kernel while an example with a conditionally positive semi-definite kernel is given below.

3.1. Semi-parametric regression with semi-positive definite kernels. Let $(\mathcal{H}_1, \langle \cdot, \cdot \rangle_1)$ be a Reproducing Kernel Hilbert Space (RKHS) with kernel $k(\mathbf{x}, \mathbf{x}')$. Also, let \mathcal{H}_0 be a Hilbert space of dimension $p < \infty$ with a basis³ given by $p_j(\mathbf{x})$ for $1 \leq j \leq p$ and such that $\mathcal{H}_0 \cap \mathcal{H}_1 = \{0\}$.

The function space that we consider is the direct sum $\mathcal{H}_0 \oplus \mathcal{H}_1$. By construction, every $f \in \mathcal{H}_0 \oplus \mathcal{H}_1$ can be decomposed uniquely as $f = f_0 + f_1$ with $f_0 \in \mathcal{H}_0$ and $f_1 \in \mathcal{H}_1$. Hence, we define the penalty functional

$$J(f) = \langle f_1, f_1 \rangle_1,$$

so that its null space is naturally $\mathcal{N}_J \triangleq \{f \in \mathcal{H}_0 \oplus \mathcal{H}_1 : J(f) = 0\} = \mathcal{H}_0$. The penalized least-squares (PLS) problem reads

$$(PLS) \quad \min_{f \in \mathcal{N}_J \oplus \mathcal{H}_1} \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 + \gamma J(f).$$

When $\mathcal{N}_J = \{0\}$, (PLS) reduces to Kernel Ridge Regression (KRR). By a classical argument⁴, the solutions of (PLS) are in the semi-parametric form $f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i) + \sum_{m=1}^p \beta_m p_m(\mathbf{x})$, where the first term includes only a *finite* number of terms. By plugging the above expression into the minimization problem (PLS), we find the discrete minimization problem

$$(3.1) \quad \min_{\boldsymbol{\alpha}, \boldsymbol{\beta}} \frac{1}{n} \|\mathbf{y} - V\boldsymbol{\beta} - K\boldsymbol{\alpha}\|_2^2 + \gamma \boldsymbol{\alpha}^\top K \boldsymbol{\alpha},$$

with $V_{im} = [p_m(\mathbf{x}_i)]$ for $1 \leq i \leq n$ and $1 \leq m \leq p$. Notice that $\frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2$ is strictly convex on \mathcal{N}_J if V is full column rank. Furthermore, (PLS) is strictly convex on $\mathcal{H}_0 \oplus \mathcal{H}_1$ if V is full column rank. Therefore, we assume in what follows that the data set is ‘unisolvent’ with respect to $(p_j)_j$, i.e., such that V is full column rank. In this case, the solution of (PLS) is unique in the light of Theorem 6.1 in Appendix (see [35]). The first order optimality condition of the optimization problem (3.1) reads

$$\begin{aligned} K[(K + n\gamma\mathbb{I})\boldsymbol{\alpha} + V\boldsymbol{\beta} - \mathbf{y}] &= 0 \\ V^\top(K\boldsymbol{\alpha} + V\boldsymbol{\beta} - \mathbf{y}) &= 0, \end{aligned}$$

where K is positive semi-definite. We assume first that K is non-singular. As it can be verified by a simple substitution of $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ in the first order conditions, the unique solution of (3.1) is obtained by solving

$$(3.2) \quad \begin{pmatrix} K + n\gamma\mathbb{I} & V \\ V^\top & 0 \end{pmatrix} \begin{pmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{pmatrix} = \begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix}.$$

Second, if K is singular, then (3.1) can have several solutions $\boldsymbol{\alpha}^*$ and $\boldsymbol{\beta}^*$, which yield the same in-sample estimator $\hat{z} = K\boldsymbol{\alpha}^* + V\boldsymbol{\beta}^*$ of the true function values $z_i = f(\mathbf{x}_i)$ for $1 \leq i \leq n$. In that case, we select the solution corresponding to the coefficients obtained by solving (3.2).

3.2. Case study: the Gaussian kernel RKHS does not contain polynomials. Let $X \subset \mathbb{R}^d$ be any set with non-empty interior and let \mathcal{H}_1 be the RKHS of the Gaussian kernel $k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|_2^2/\sigma^2)$ defined on $X \times X$. Under these assumptions, Theorem 2 in [44] states that \mathcal{H}_1 does not contain any polynomial on X , including the non-zero constant function. We can then solve the functional minimization problem (PLS) where \mathcal{H}_0 is a finite set of polynomials, since $\mathcal{H}_0 \cap \mathcal{H}_1 = \{0\}$.

EXAMPLE 2 (LS-SVM with the Gaussian kernel). *In particular, the constant $p_1(\mathbf{x}) = 1$ is not part of the RKHS of the Gaussian kernel. Therefore, Least-Squares Support Vector Machine (LS-SVM) [56] with the Gaussian kernel is also a particular case of the above discussion. Its dual optimization problem indeed reads $\min_{\boldsymbol{\alpha}, b} \frac{1}{n} \|\mathbf{y} - K\boldsymbol{\alpha} - b\mathbf{1}_n\|_2^2 + \gamma \boldsymbol{\alpha}^\top K \boldsymbol{\alpha}$, where the real b is the so-called bias term.*

³Any finite dimensional vector space can be endowed with a suitable scalar product so that is also a RKHS. Let $(\mathcal{H}_0, \langle \cdot, \cdot \rangle_0)$ be a RKHS. Then, \mathcal{H}_0 is the orthogonal complement of \mathcal{H}_1 with respect to the following inner product: $\langle f_0 + g_0, f_1 + g_1 \rangle \triangleq \langle f_0, g_0 \rangle_0 + \langle f_1, g_1 \rangle_1$.

⁴A ‘representer theorem’, see, e.g. [35, Section 2.3.2] or [48].

In Figure 1, we illustrate the use of semi-parametric regression with a Gaussian kernel on a toy example consisting of a linear trend with two Gaussian bumps, i.e., $f(x) = x + 7 + 4 \exp(-(x-4)^2) - 4 \exp(-(x+4)^2)$. The training points are sampled uniformly within the interval $[-10, 10]$ and the function samples are $y_i = f(x_i) + \epsilon_i$ where $\epsilon_i \sim \mathcal{N}(0, 0.2)$ for $1 \leq i \leq n$ and $n = 40$. The test set consists of 1000 points sampled uniformly in the interval $[-11, 11]$. This construction allows to assess the ability of the estimated function to capture the linear trend of the ground truth.

Let \mathcal{H}_1 be the RKHS of the Gaussian kernel. We compare the results obtained by different choices of the space of polynomials \mathcal{H}_0 . Specifically, the following models are estimated: **Model 1**: $\hat{f}(x) = \beta_1 + \beta_2 x + \sum_i \alpha_i k(x, x_i)$ (semi-parametric), **Model 2**: $\hat{f}(x) = \beta_1 + \sum_i \alpha_i k(x, x_i)$ (LS-SVM) and **Model 3**: $\hat{f}(x) = \sum_i \alpha_i k(x, x_i)$ (KRR). For all the models drawn in Figure 1a, the bandwidth is fixed to $\sigma = 1$ to match the width of the two Gaussians of the ground truth and the regularization parameter γ takes values in the set $\{10^{-j}\}_{j \in \{1, \dots, 8\}}$. For completeness, we also include the performance of each model after parameter tuning in Figure 1b (dashed line), where both the bandwidth $\sigma \in \{0.1, 0.2, \dots, 0.9, 1, 2, 3, \dots, 10\}$ and regularization parameter are determined by using 10 fold cross-validation, on the above-mentioned grid. The simulation is repeated 25 times and the error bars show the 97.5% confidence interval. In Figure 1a, the function prediction in low density regions of both **model 2** and **model 3** quickly moves to the bias. This is avoided by including a parametric linear part such as in **model 1**. We emphasize that the differences between the three models are reduced within the interval $[-10, 10]$ if the number of training samples becomes larger.

3.3. Case study: conditionally positive semi-definite kernels and thin-plate splines. A well-known choice of penalty functional yielding to a linear system with a conditionally positive semi-definite kernel is associated to thin-plate splines, and is given by

$$J_p^d(f) \triangleq \langle f, f \rangle_1 \text{ with } \langle f, g \rangle_1 = \sum_{\alpha_1 + \dots + \alpha_d = m} \frac{p!}{\alpha_1! \dots \alpha_d!} \int_{\mathbb{R}^d} \frac{\partial^p f}{\partial x_1^{\alpha_1} \dots \partial x_d^{\alpha_d}} \frac{\partial^p g}{\partial x_1^{\alpha_1} \dots \partial x_d^{\alpha_d}} dx_1 \dots dx_d,$$

where $J_p^d(f)$ is a squared semi-norm on $\{f : J_p^d(f) < \infty\}$ for a large enough regularity index with respect to the dimension, i.e., for $2p > d$. Its null space \mathcal{N}_J consists of polynomials of maximal total order equal to $p - 1$. Following section 4.3.2 of [35], the penalty functional $J_p^d(f)$ satisfies

$$J_p^d \left(\sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i) \right) = \sum_{i,j=1}^n \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j), \text{ for all } \boldsymbol{\alpha} \in \mathbb{R}^n \text{ such that } V^\top \boldsymbol{\alpha} = 0,$$

where V is assumed to be full column rank and where the thin-plate spline kernel reads

$$k(\mathbf{x}, \mathbf{x}') = \begin{cases} \|\mathbf{x} - \mathbf{x}'\|_2^{2p-d} \log \|\mathbf{x} - \mathbf{x}'\|_2 & \text{for even } d \\ \|\mathbf{x} - \mathbf{x}'\|_2^{2p-d} & \text{for odd } d. \end{cases}$$

The kernel $k(\mathbf{x}, \mathbf{x}')$ is conditionally positive semi-definite, namely, $\sum_{i,j} \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) \geq 0$, for all vectors satisfying $V^\top \boldsymbol{\alpha} = 0$. Again, V is assumed here to be full column rank. By an similar argument as in the previous section (see [35]), the solution of the least-squares penalized regression

$$\min_{J_p^d(f) < \infty} \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 + \gamma J_p^d(f),$$

is of the form $f(\mathbf{x}) = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i) + \sum_{j=1}^p \beta_j p_j(\mathbf{x})$ with $V^\top \boldsymbol{\alpha} = 0$. This result is proved in [25, Theorem 4 bis] in the case of optimal interpolation (i.e., $\gamma \rightarrow 0$). The substitution of $f(\mathbf{x})$ into the objective above yields a similar discrete minimization problem as (3.1) with the extra condition $V^\top \boldsymbol{\alpha} = 0$. In particular, a solution of this minimization problem is given by the same system as (3.2) with the exception that, here, K is *conditionally* positive semi-definite. Let V_\perp a matrix with orthonormal columns such that $\mathbb{P}_{V_\perp} = V_\perp V_\perp^\top$. The solution of this linear system is given as

$$\begin{aligned} \boldsymbol{\alpha}^* &= V_\perp (V_\perp^\top K V_\perp + n\gamma \mathbb{I}_{n-p})^{-1} V_\perp^\top \mathbf{y} \\ \boldsymbol{\beta}^* &= (V^\top V)^{-1} V^\top (\mathbf{y} - K \boldsymbol{\alpha}^*), \end{aligned}$$

where we used the fact that V is full column rank and $\mathbb{P}_{V^\perp} K \mathbb{P}_{V^\perp} \succeq 0$. Notice that the full in-sample estimator is $\hat{\mathbf{z}} = \tilde{K}(\tilde{K} + n\gamma\mathbb{I})^{-1}\mathbf{y} + \mathbb{P}_V\mathbf{y}$. The RKHS associated to the thin-plate splines and a *psd* kernel built from the conditionally positive semi-definite kernel are determined in [35], where this problem is put in the form of (PLS). Then, the domain $\{f : J_p^d(f) < \infty\}$ is shown to be the direct sum of a RKHS and the space of polynomials of maximal total degree $p - 1$. For completeness, let us mention that a discussion of optimal interpolation with conditionally positive semi-definite kernels, within the framework of Hilbertian subspaces of L. Schwartz, can be found in the PhD thesis [32]. Consider now the use of extended L -ensembles for obtaining designs in the context of optimal interpolations.

4. Implicit regularization of optimal interpolation with a determinantal design. In the context of the applications mentioned in the introduction, we consider here the problem of interpolating function values given a small training data set. In this section, we assume that obtaining the responses y_i is expensive and therefore we look for a discrete design $(\mathbf{x}_{i_\ell}, y_{i_\ell})$ for $i_\ell \in \mathcal{C}$. An interpolator is obtained by taking the ‘ridgeless’ limit $\gamma \rightarrow 0$ of the regression problem (3.1). Let $\mathbf{k}_x = [k(\mathbf{x}, \mathbf{x}_1) \dots k(\mathbf{x}, \mathbf{x}_n)]^\top \in \mathbb{R}^n$ and $\mathbf{p}_x = [p_1(\mathbf{x}) \dots p_p(\mathbf{x})]^\top \in \mathbb{R}^p$ for all $\mathbf{x} \in \mathbb{R}^d$. The estimated interpolator on a subset \mathcal{C} reads

$$(4.1) \quad \hat{f}_0(\mathbf{x}, \mathcal{C}) = (\mathbf{k}_x^\top \quad \mathbf{p}_x^\top) \begin{pmatrix} K_{\mathcal{C}\mathcal{C}} & V_{\mathcal{C}} \\ V_{\mathcal{C}}^\top & 0 \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{y}_{\mathcal{C}} \\ \mathbf{0} \end{pmatrix}.$$

If $\mathcal{C} \sim DPP(K, V)$, notice that: (i) the square matrix on the RHS of (4.1) is non-singular almost surely, and (ii) $V_{\mathcal{C}}$ is full column rank almost surely, in the light of Lemma 6.10. One of our main results is that the linear system in (3.2) is regularized on expectation when the subsets \mathcal{C} are sampled according to a suitable DPP.

THEOREM 4.1 (implicit regularization on expectation). *Let (K, V) be a NNP. Let $\mathbf{u}_0, \mathbf{v}_0 \in \mathbb{R}^n$ and $\mathbf{u}_1, \mathbf{v}_1 \in \mathbb{R}^p$. Then, we have the following identity*

$$\mathbb{E}_{\mathcal{C} \sim DPP(K, V)} \left[\begin{pmatrix} \mathbf{u}_{0, \mathcal{C}} \\ \mathbf{u}_1 \end{pmatrix}^\top \begin{pmatrix} K_{\mathcal{C}\mathcal{C}} & V_{\mathcal{C}} \\ V_{\mathcal{C}}^\top & 0 \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{v}_{0, \mathcal{C}} \\ \mathbf{v}_1 \end{pmatrix} \right] = \begin{pmatrix} \mathbf{u}_0 \\ \mathbf{u}_1 \end{pmatrix}^\top \begin{pmatrix} K + \mathbb{I} & V \\ V^\top & 0 \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{v}_0 \\ \mathbf{v}_1 \end{pmatrix}.$$

Proof. The proof of this result is given in Section 6.3 in Appendix and mainly relies on the matrix determinant lemma. \square

Notice that only the upper left block is regularized in the above matrix inverse. Also, this identity remains valid when K is replaced by $\tilde{K} = \mathbb{P}_{V^\perp} K \mathbb{P}_{V^\perp}$. Similarly to (4.1), the γ -regularized regressor on the full data set obtained by solving (PLS) is

$$\hat{f}_\gamma(\mathbf{x}) = (\mathbf{k}_x^\top \quad \mathbf{p}_x^\top) \begin{pmatrix} K + n\gamma\mathbb{I} & V \\ V^\top & 0 \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{y} \\ \mathbf{0} \end{pmatrix}.$$

The upshot is that the interpolator obtained with this determinantal design is actually regularized on expectation, as a direct consequence of Theorem 4.1. This result is formalized in Corollary 4.2 which generalizes a similar result for KRR in the ridgeless limit given in [54].

COROLLARY 4.2 (ensemble of interpolators). *Let $\mathcal{C} \sim DPP(K/(n\gamma), V)$. We have $\mathbb{E}_{\mathcal{C}}[\hat{f}_0(\mathbf{x}, \mathcal{C})] = \hat{f}_\gamma(\mathbf{x})$ for all $\mathbf{x} \in \mathbb{R}^d$.*

The interpretation of Corollary 4.2 goes as follows: an average of interpolators obtained with this random design gives a regularized regressor on the full data set. We refer to [17] for another discussion of connections between Bayesian experimental design and DPPs. In the next section, we illustrate the use of a discrete determinantal design in the context of optimal interpolation with thin-plate splines.

4.1. Illustration of a discrete determinantal design for thin-plate spline interpolation. We illustrate the effect of subsampling for thin-plate spline interpolation on Franke’s function, which is frequently used to demonstrate radial basis function interpolation problems. Franke’s function has two Gaussian peaks

of different heights, and a smaller dip:

$$f(\mathbf{x}) = 0.75 \exp\left(-\frac{(9x_1 - 2)^2}{4} - \frac{(9x_2 - 2)^2}{4}\right) + 0.75 \exp\left(-\frac{(9x_1 + 1)^2}{49} - \frac{9x_2 + 1}{10}\right) \\ + 0.5 \exp\left(-\frac{(9x_1 - 7)^2}{4} - \frac{(9x_2 - 3)^2}{4}\right) - 0.2 \exp\left(- (9x_1 - 4)^2 - (9x_2 - 7)^2\right).$$

The full training set consists of 5000 points sampled uniformly at random within $[0, 1]^2$, the test set consists of 10000 points sampled from the same domain. The full interpolation problem is solved by using (3.2) with $\gamma = 0$ and regression function $f(\mathbf{x}) = \sum_{m=1}^2 b_m x_m + b_0 + \sum_{i=1}^n \alpha_i \|\mathbf{x} - \mathbf{x}_i\|_2^2 \log \|\mathbf{x} - \mathbf{x}_i\|_2$. The subsampled interpolation problem is solved by using (4.1). In this simulation, we compare uniform sampling to the associated partial-projection DPP. The simulation is repeated for an increasing subset size, and the performance is measured by the mean squared error (MSE) on the test set. In the experiments, we sample each time a fixed size partial-projection DPP (with $|\mathcal{C}| = k$) to finer control the number of sampled landmarks. Every sampling is repeated 10 times and the averaged results are visualized in Figure 2. Error bars correspond to the 97.5% confidence interval. For a given subset size, the partial-projection DPP outperforms uniform sampling.

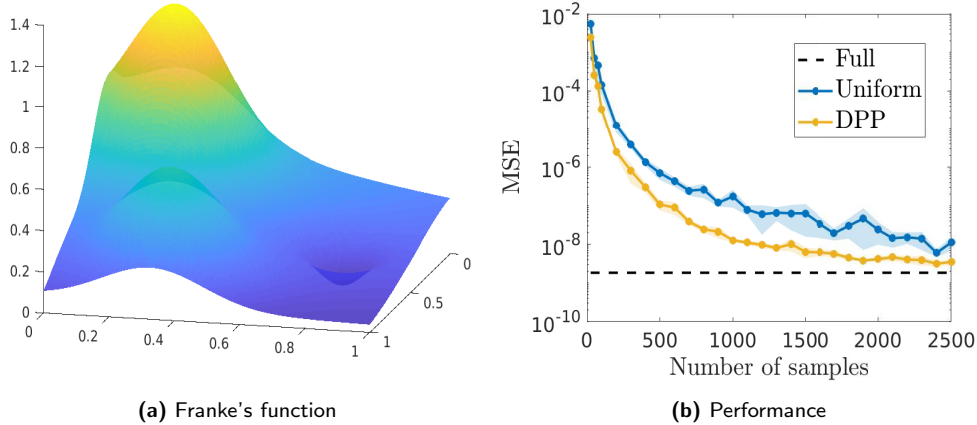


Fig. 2: Figure 2a displays a mesh plot of the Franke's function. The MSE on the test set as a function of the subset size $|\mathcal{C}|$ is given in Figure 2b.

4.2. Empirical results for Gaussian kernel interpolation. We illustrate here the effect of extended L -ensemble sampling versus uniform sampling for subsampled interpolation on a number of UCI benchmark regression data sets: **Boston Housing**, **Abalone** and **Parkinson**. Both the regressors and response are standardized, afterwards the data set is split into a 50% training and 50% test set, and the performance is measured by the total MSE: $\sum_{i=1}^n \|y_i - \hat{z}_i\|^2$. To obtain the regressor, we solve the system (3.2). A Gaussian kernel $k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|_2^2 / \sigma^2)$ is used with a linear regression component: $V = [X \mathbf{1}_n]$ where $X = [\mathbf{x}_1 \dots \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$. The squared bandwidth is determined by using the median heuristic [34], computed as: $\hat{\sigma}^2 = \text{median}\{\|x_i - x_j\|_2^2 : 1 \leq i < j \leq n\} / 2$. Cross-validation is not possible as the full \mathbf{y} is hidden from us in the experimental design setup. The simulation is repeated 25 times and the error bars show the 97.5% confidence interval. The results are displayed in Figure 3. We observe that extended L -ensemble sampling improves the performance especially for smaller number of samples, compared to uniform sampling.

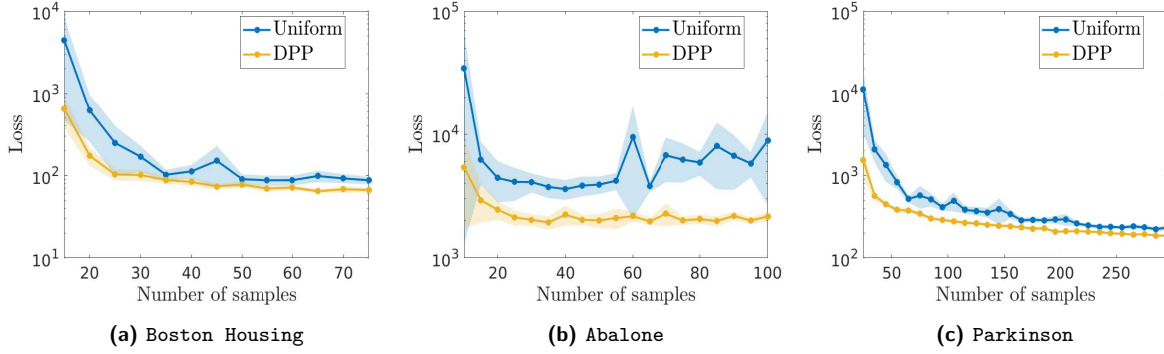


Fig. 3: The total loss (MSE) in function of the number of landmarks using uniform vs extended DPP sampling with the bandwidth estimated using the median heuristic.

5. Large scale regularized semi-parametric regression. In this section, we consider the setting where the training points (\mathbf{x}_i, y_i) with $1 \leq i \leq n$ are abundant. Penalized least-squares problems of the form of (PLS) have solutions which are determined by $n + p$ parameters where n is the number of training points and p is the number of functions used in the parametric component. As it was anticipated in the introduction, when n is large, it can be interesting to reduce the number of parameters describing the estimated function, for instance, in order to allow for a faster out-of-sample prediction.

5.1. Preliminary results about implicit regularization. For a $n \times n$ positive semi-definite kernel matrix K , the Nyström approximation of K associated to $\mathcal{C} \subseteq [n]$ reads $K_{\mathcal{C}}^{\top} K_{\mathcal{C}\mathcal{C}}^{+} K_{\mathcal{C}}$. We extend here this definition to the case of conditionally positive semi-definite kernels. First, we provide in Proposition 5.1 the expectation of an analogue of $CK_{\mathcal{C}\mathcal{C}}^{+}C^{\top}$ accounting for conditional positivity.

PROPOSITION 5.1 (implicit regularization of the projected kernel matrix). *Let $\mathcal{C} \sim DPP(K/\lambda, V)$. Then, we have*

$$(5.1) \quad \mathbb{E}_{\mathcal{C}}[I(\mathcal{C})] = (\tilde{K} + \lambda \mathbb{P}_{V^{\perp}})^{+} \text{ with } I(\mathcal{C}) = (C\mathbb{P}_{V_{\mathcal{C}}^{\perp}}K_{\mathcal{C}\mathcal{C}}\mathbb{P}_{V_{\mathcal{C}}^{\perp}}C^{\top})^{+}.$$

Furthermore, it also holds that

$$(5.2) \quad \mathbb{E}_{\mathcal{C}}[(CV_{\mathcal{C}})^{+}(\mathbb{I} - KI(\mathcal{C}))] = V^{+}(\mathbb{I} - K(\tilde{K} + \lambda \mathbb{P}_{V^{\perp}})^{+})$$

$$(5.3) \quad \mathbb{E}_{\mathcal{C}}[(CV_{\mathcal{C}})^{+}(K - KI(\mathcal{C})K)(CV_{\mathcal{C}})^{+\top}] = V^{+}(K + \lambda \mathbb{I} - K(\tilde{K} + \lambda \mathbb{P}_{V^{\perp}})^{+}K)V^{+\top},$$

The identity (5.1) is equivalent to the expected pseudo-inverse formula (1.4) announced in the introduction whereas the identities (5.2) and (5.3) are incidental and are studied in more detail the the next section.

Proof. We begin by noticing that, thanks to Lemma 6.7 in Appendix, it holds that

$$(C\mathbb{P}_{V_{\mathcal{C}}^{\perp}}K_{\mathcal{C}\mathcal{C}}\mathbb{P}_{V_{\mathcal{C}}^{\perp}}C^{\top})^{+} = C(\mathbb{P}_{V_{\mathcal{C}}^{\perp}}K_{\mathcal{C}\mathcal{C}}\mathbb{P}_{V_{\mathcal{C}}^{\perp}})^{+}C^{\top} \text{ and } (\tilde{K} + \lambda \mathbb{P}_{V^{\perp}})^{+} = \mathbb{P}_{V^{\perp}}(\tilde{K} + \lambda \mathbb{I})^{-1}\mathbb{P}_{V^{\perp}}.$$

Without loss of generality, we now take $\lambda = 1$, since the results can be recovered at the end for any $\lambda > 0$ by a simple rescaling of K . Let $B(\mathcal{C})$ be a matrix whose columns are an orthonormal basis of the column space of $(V_{\mathcal{C}})^{\perp}$. In this case, we have $\mathbb{P}_{V_{\mathcal{C}}^{\perp}} = B(\mathcal{C})B(\mathcal{C})^{\top}$. Then, we have the explicit expression

$$(\mathbb{P}_{V_{\mathcal{C}}^{\perp}}K_{\mathcal{C}\mathcal{C}}\mathbb{P}_{V_{\mathcal{C}}^{\perp}})^{+} = B(\mathcal{C})(B(\mathcal{C})^{\top}K_{\mathcal{C}\mathcal{C}}B(\mathcal{C}))^{-1}B(\mathcal{C}),$$

where we used once more Lemma 6.7 (with $S = B(\mathcal{C})$ and $M = B(\mathcal{C})^{\top}K_{\mathcal{C}\mathcal{C}}B(\mathcal{C})$). The identities (5.1), (5.2) and (5.3) are obtained in what follows by merely calculating matrix inverse in the formula given in Theorem 4.1. For simplicity, define

$$T(\mathcal{C}) = \begin{pmatrix} C & 0 \\ 0 & \mathbb{I} \end{pmatrix} \begin{pmatrix} K_{\mathcal{C}\mathcal{C}} & V_{\mathcal{C}} \\ V_{\mathcal{C}}^{\top} & 0 \end{pmatrix}^{-1} \begin{pmatrix} C^{\top} & 0 \\ 0 & \mathbb{I} \end{pmatrix}, \text{ for } \mathcal{C} \sim DPP(K, V).$$

The above matrix inverse is now calculated by using Lemma 6.6 in Appendix (with $A = K_{CC}$ and $W = V_C$). Remark that $B^\top(\mathcal{C})K_{CC}B(\mathcal{C})$ is non-singular. Indeed, since $\mathcal{C} \sim DPP(K, V)$, we have

$$0 \neq \det \begin{pmatrix} K_{CC} & V_C \\ V_C^\top & 0 \end{pmatrix} = (-1)^p \det(V_C^\top V_C) \underbrace{\det(B^\top(\mathcal{C})K_{CC}B(\mathcal{C}))}_{\neq 0},$$

where the determinant on the LHS is calculated thanks to Lemma 6.10 in Appendix. Thus, we find the following expression

$$T(\mathcal{C}) = \begin{pmatrix} C(\mathbb{P}_{V_C^\perp} K_{CC} \mathbb{P}_{V_C^\perp})^+ C^\top & C(\mathbb{I} - (\mathbb{P}_{V_C^\perp} K_{CC} \mathbb{P}_{V_C^\perp})^+ K_{CC}) V_C^{+\top} \\ V_C^+(\mathbb{I} - K_{CC}(\mathbb{P}_{V_C^\perp} K_{CC} \mathbb{P}_{V_C^\perp})^+) C^\top & -V_C^\top (K_{CC} - K_{CC}(\mathbb{P}_{V_C^\perp} K_{CC} \mathbb{P}_{V_C^\perp})^+ K_{CC}) V_C^{+\top} \end{pmatrix}.$$

Next, by using Theorem 4.1, it holds that $\mathbb{E}_{\mathcal{C}}[T(\mathcal{C})] = (K + \mathbb{I} \ V; V^\top \ 0)^{-1}$, where the inverse matrix can be calculated thanks to Lemma 6.6, as follows:

$$\begin{pmatrix} K + \mathbb{I} & V \\ V^\top & 0 \end{pmatrix}^{-1} = \begin{pmatrix} (\mathbb{P}_{V^\perp}(K + \mathbb{I})\mathbb{P}_{V^\perp})^+ & (\mathbb{I} - (\mathbb{P}_{V^\perp}(K + \mathbb{I})\mathbb{P}_{V^\perp})^+ K)V^{+\top} \\ V^+(\mathbb{I} - K(\mathbb{P}_{V^\perp}(K + \mathbb{I})\mathbb{P}_{V^\perp})^+) & -V^+(K + \mathbb{I} - K(\mathbb{P}_{V^\perp}(K + \mathbb{I})\mathbb{P}_{V^\perp})^+ K)V^{+\top} \end{pmatrix}.$$

The desired result follows by identifying the different terms in the above block matrix, and by using the simply identity $(\mathbb{P}_{V^\perp}(K + \mathbb{I})\mathbb{P}_{V^\perp})^+ = (\tilde{K} + \mathbb{P}_{V^\perp})^+$. This completes the proof. \square

The above cumbersome expressions in (5.2) and (5.3) do not seem at first sight to have a straightforward interpretation. Nonetheless, for a special choice of partial projection DPP, the identities in (5.2) and (5.3) can be used in the context of random design regression in the same spirit as rescaled volume sampling [20], as we discuss in the following interlude subsection.

5.2. Interlude: an unbiased estimator for random design regression with a DPP. In order to illustrate the interest of (5.2) and (5.3) in Proposition 5.1, we consider the connection between extended L -ensembles, volume sampling [22] and projection DPPs (see e.g., [6]). This yields a partial projection DPP extending the well-known volume sampling method, which has often been considered in the literature for example in randomized linear algebra [1]. Another related approach called ‘proportional volume sampling’ has been discussed in [49] for general Polish spaces, while a related Bayesian approach is given in [17]. The generalization presented here can be used in order to find discrete random designs for parametric regression problems, as we explain below.

A random-size volume sampling. First, we recall the definition of volume sampling.

DEFINITION 5.2 (volume sampling). *Let $V \in \mathbb{R}^{n \times p}$ be a matrix with full column rank. The probability to volume sample a subset $\mathcal{C}_0 \sim \text{Vol}_k(V)$ of size $k \geq p$ is*

$$\Pr(\mathcal{C}_0) = \frac{\det(V_{\mathcal{C}_0}^\top V_{\mathcal{C}_0})}{\binom{n-p}{k-p} \det(V^\top V)}.$$

For $k = p$, volume sampling is a projection DPP with marginal kernel $V(V^\top V)^{-1}V^\top$.

Next, consider the partial projection $DPP(t\mathbb{I}, V)$ for $t > 0$ corresponding to the marginal kernel

$$P = q\mathbb{P}_{V^\perp} + \mathbb{P}_V \text{ where } q = \frac{t}{1+t}.$$

This process is in fact a rescaling of volume sampling, that is,

$$(5.4) \quad \Pr(\mathcal{Y} = \mathcal{C}) = \frac{\det(t\mathbb{I}_{|\mathcal{C}|})}{\det(\mathbb{P}_V + (1+t)\mathbb{P}_{V^\perp})} \frac{\det(V_{\mathcal{C}}^\top V_{\mathcal{C}})}{\det(V^\top V)} = q^{|\mathcal{C}|-p} (1-q)^{n-|\mathcal{C}|} \times \frac{\det(V_{\mathcal{C}}^\top V_{\mathcal{C}})}{\det(V^\top V)},$$

where the first equality uses the formula (2.1) for the probability mass function of \mathcal{C} . There is a clear analogy with the volume sampling method of [22, Eq. (1)], although the sample size of (5.4) is here a random variable which satisfies

$$(5.5) \quad \frac{\mathbb{E}_{\mathcal{C}}[|\mathcal{C}|] - p}{n - p} = \frac{t}{1+t}, \text{ for } \mathcal{C} \sim DPP(t\mathbb{I}, V)$$

while volume sampling always returns subsets of a fixed size. In the same spirit as in [22], we propose below a sampling strategy based on a combination of volume sampling and a Bernoulli process, and whose correctness is discussed in Proposition 5.4.

DEFINITION 5.3 (Bernoulli process). *Let $0 \leq q \leq 1$ and m be a positive integer. A Bernoulli(q) process over $\{0, 1\}^m$ is a sequence of i.i.d. Bernoulli random variables with success probability q , so that the probability to observe any sequence (s_1, \dots, s_m) is $q^k(1-q)^{m-k}$ where k is the number of successes in the sequence.*

Commonly, we define a Bernoulli process over an ordered sequence of m items by selecting an item of its corresponding entry in $(s_1, \dots, s_m) \in \{0, 1\}^m$ is equal to unity. A finite Bernoulli(q) process is a DPP with marginal kernel $P = q\mathbb{I}$ for $0 \leq q \leq 1$. Naturally, the expected sample size of Bernoulli(q) over a set of m items is mq .

PROPOSITION 5.4 (two-stage sampling method for (5.4)). *A sample distributed according to $DPP(t\mathbb{I}, V)$ can be obtained as follows:*

1. draw a volume sample $\mathcal{C}_0 \sim \text{Vol}_p(V)$ of the set $[n]$,
2. draw a sample $\mathcal{R} \subseteq [n] \setminus \mathcal{C}_0$ according to Bernoulli($\frac{t}{t+1}$),

and return $\mathcal{C} = \mathcal{C}_0 \cup \mathcal{R}$.

Proof. Let $\mathcal{C} \subseteq [n]$ such that $|\mathcal{C}| \geq p$. Consider all the decompositions $\mathcal{C} = \mathcal{R} \cup \mathcal{C}_0$ where $\mathcal{R} = \mathcal{C} \setminus \mathcal{C}_0$, for all the subsets \mathcal{C}_0 such that $|\mathcal{C}_0| = p$. Notice that $|\mathcal{R}| = |\mathcal{C}| - p$. Then, the probability that \mathcal{C} is obtained by combining a volume sample \mathcal{C}_0 with a sample \mathcal{R} drawn by a Bernoulli process on $[n] \setminus \mathcal{C}_0$ is

$$\Pr(\mathcal{C}) = \sum_{\mathcal{C}_0 \subseteq \mathcal{C}: |\mathcal{C}_0|=p} \underbrace{q^{|\mathcal{C}|-p} (1-q)^{n-|\mathcal{C}|}}_{\Pr(\mathcal{R}=\mathcal{C} \setminus \mathcal{C}_0 | \mathcal{C}_0)} \underbrace{\det(V_{\mathcal{C}_0}^\top V_{\mathcal{C}_0}) / \det(V^\top V)}_{\Pr(\mathcal{C}_0) \text{ (volume sampling)}} = q^{|\mathcal{C}|-p} (1-q)^{n-|\mathcal{C}|} \det(V_{\mathcal{C}}^\top V_{\mathcal{C}}),$$

where the second equality follows from the Cauchy-Binet identity $\sum_{\mathcal{C}_0 \subseteq \mathcal{C}: |\mathcal{C}_0|=p} \det(V_{\mathcal{C}_0}^\top V_{\mathcal{C}_0}) = \det(V_{\mathcal{C}}^\top V_{\mathcal{C}})$ (see, Lemma 6.4 in Appendix). \square

The following remarks are in order. For $\mathcal{C} \sim DPP(t\mathbb{I}, V)$, Proposition 5.1 reduces to the simple identity $\mathbb{E}_{\mathcal{C}}[C\mathbb{P}_{V_{\mathcal{C}}^+} C^\top] = \frac{t}{1+t}\mathbb{P}_{V^\perp}$. By using the expression of the marginal kernel of $DPP(t\mathbb{I}, V)$, an elementary manipulation yields the following expression

$$(5.6) \quad \mathbb{E}_{\mathcal{C}}[CV_{\mathcal{C}}(CV_{\mathcal{C}})^+] = VV^+,$$

where $(CV_{\mathcal{C}})^+ = V_{\mathcal{C}}^+ C^\top$. This can be interpreted as follows: the expectation of the projector onto the column space of $CV_{\mathcal{C}}$ is the projector onto the column space of V . An identity analogous to (5.6) has been obtained in the framework of matroids in [38] which revisited a related identity in [42, Thm 1]. In the same spirit as in (5.6), the identities in Proposition 5.1 simplify in the special case $\mathcal{C} \sim DPP(t\mathbb{I}, V)$ to

$$(5.7) \quad \mathbb{E}_{\mathcal{C}}[(CV_{\mathcal{C}})^+] = V^+, \text{ and } \mathbb{E}_{\mathcal{C}}[\underbrace{(CV_{\mathcal{C}})^+(CV_{\mathcal{C}})^{\top}}_{(V_{\mathcal{C}}^\top V_{\mathcal{C}})^{-1}}] = \frac{n-p}{\mathbb{E}_{\mathcal{C}}[|\mathcal{C}|] - p} \times \frac{V^+ V^{\top}}{(V^\top V)^{-1}},$$

where we used the expression relating $t > 0$ with the expected subset size (5.5). The above expressions in (5.7) can be used to define an unbiased estimator for subsampled regression problems over a smaller set of inputs and outputs, as we explain below.

Discrete random design regression. Let $\mathbf{y} \in \mathbb{R}^n$ and consider the (full-fledged) parametric regression problem

$$\boldsymbol{\beta}^* = \arg \min_{\boldsymbol{\beta}} \|V\boldsymbol{\beta} - \mathbf{y}\|_2^2,$$

with $V \in \mathbb{R}^{n \times p}$, and whose solution writes $\mathbf{w}^* = V^+ \mathbf{y}$. It is known that solving a smaller regression problem with sampled input-output pairs $(V_{\mathcal{C}}, \mathbf{y}_{\mathcal{C}})$ regardless of the outputs for $\mathcal{C} \subseteq [n]$ can lead to a biased estimator of the full regressor [21], for instance if \mathcal{C} is sampled uniformly at random. However, for $\mathcal{C} \sim DPP(t\mathbb{I}, V)$ the estimator $\hat{\boldsymbol{\beta}}(\mathcal{C}) = V_{\mathcal{C}}^+ \mathbf{y}_{\mathcal{C}}$ of the quantity $V^+ \mathbf{y}$ is unbiased, while the matrix variance writes

$$(5.8) \quad \mathbb{E}[(CV_{\mathcal{C}})^+(CV_{\mathcal{C}})^{\top}] - \mathbb{E}[(CV_{\mathcal{C}})^+] \mathbb{E}[(CV_{\mathcal{C}})^+]^{\top} = \frac{n - \mathbb{E}[|\mathcal{C}|]}{\mathbb{E}[|\mathcal{C}|] - p} \times V^+ V^{\top},$$

thanks to (5.7). A main difference with the designs of [21, 20] obtained with a fixed-size rescaled volume sampling is that here: (i) the subset size is random, and (ii) the above variance grows unbounded as $t \rightarrow 0$. On the contrary, a large value of t promotes a large expected subset size and a small variance.

Let us recall some related results for volume sampling which yield formulae analogous to (5.7). On the one hand, the identity $\mathbb{E}[(CV_C)^+] = V^+$ for $\mathcal{C} \sim \text{Vol}_k(V)$ has been obtained in the context of random design regression in [21] and in [20, Thm 2.10]. On the other hand, the result of [23, Eq (1)] yields the ‘second moment’ $\mathbb{E}[(CV_C)^+(CV_C)^{+\top}] = \frac{n-p+1}{k-p+1}V^+V^{+\top}$ for $\mathcal{C} \sim \text{Vol}_k(V)$ which has a slightly different form in comparison with our result in (5.7). Indeed, volume sampling for $k = p$ yields a finite variance estimator while the variance of the estimator built with our partial projection DPP grows unbounded as the expected subset size goes to p .

As it is discussed in [49], it seems that random design methods only outperform existing optimal design methods [50] in very specific settings. We leave a more detailed empirical exploration of our random-size volume sampling for further work.

5.3. Projected Nyström approximation. In the light of Proposition 5.1, we now define the projected Nyström approximation as a generalization of the common Nyström approximation which was given in the introduction.

DEFINITION 5.5 (projected Nyström approximation). *Let (K, V) be a NNP and $\mathcal{C} \subseteq [n]$. Let $\tilde{K} = \mathbb{P}_{V^\perp} K \mathbb{P}_{V^\perp}$. The projected Nyström approximation⁵ of \tilde{K} is*

$$\widetilde{L}(\mathcal{C}) = \mathbb{P}_{V^\perp} K_{\mathcal{C}}^\top B(\mathcal{C}) \left(B^\top(\mathcal{C}) K_{\mathcal{C}\mathcal{C}} B(\mathcal{C}) \right)^+ B^\top(\mathcal{C}) K_{\mathcal{C}} \mathbb{P}_{V^\perp},$$

where $B(\mathcal{C}) \in \mathbb{R}^{k \times (k-p)}$ be a matrix whose columns are an orthonormal basis of $(V_{\mathcal{C}})^\perp$.

Several remarks are in order. First, the pseudo-inverse in the above definition can be replaced by a matrix inverse almost surely if $\mathcal{C} \sim \text{DPP}(K, V)$. Second, it is worth emphasizing that, in Definition 5.5, the submatrices $K_{\mathcal{C}}$ and $K_{\mathcal{C}\mathcal{C}}$ are sufficient to construct the projected Nyström approximation while \tilde{K} should not be explicitly constructed. Therefore, the projected Nyström approximation is also promising in order to solve problems where the kernel matrix is too large compared to the computer memory. Third, the projected Nyström approximation given in Definition 5.5 satisfies the following desirable property: the submatrices $\widetilde{L}(\mathcal{C})_{\mathcal{C}\mathcal{C}}$ and $\tilde{K}_{\mathcal{C}\mathcal{C}}$ match in the appropriate subspace, that is, $B^\top(\mathcal{C}) \widetilde{L}(\mathcal{C})_{\mathcal{C}\mathcal{C}} B(\mathcal{C}) = B^\top(\mathcal{C}) \tilde{K}_{\mathcal{C}\mathcal{C}} B(\mathcal{C})$, for $\mathcal{C} \subseteq [n]$. This is a consequence of Corollary 5.6 given below.

COROLLARY 5.6 (Nyström approximation error). *Let $\mathcal{C} \sim \text{DPP}(K/\lambda, V)$ with $\lambda > 0$. Let $B(\mathcal{C}) \in \mathbb{R}^{|\mathcal{C}| \times (|\mathcal{C}|-p)}$ be a matrix whose columns are an orthonormal basis of $(V_{\mathcal{C}})^\perp$. Then, we have the following identities*

$$(i) \quad 0 \preceq \widetilde{L}(\mathcal{C}) = \tilde{K}_{\mathcal{C}}^\top B(\mathcal{C}) \left(B^\top(\mathcal{C}) \tilde{K}_{\mathcal{C}\mathcal{C}} B(\mathcal{C}) \right)^{-1} B^\top(\mathcal{C}) \tilde{K}_{\mathcal{C}} \preceq \tilde{K},$$

$$(ii) \quad \mathbb{E}[\tilde{K} - \widetilde{L}(\mathcal{C})] = \lambda \tilde{K} (\tilde{K} + \lambda \mathbb{I})^{-1} \preceq \lambda \mathbb{I}.$$

Notice that Corollary 5.6 shows that the expected error of the approximation naturally decreases if the expected number of sampled landmarks increases, that is, as $\lambda > 0$ goes to zero. The proof of this result merely follows from Proposition 5.1.

Proof of Corollary 5.6. (i) It is easy to check that $\widetilde{L}(\mathcal{C}) \succeq 0$. The first identity is simply obtained by using Lemma 6.12 in Appendix. To show that $\widetilde{L}(\mathcal{C}) \preceq \tilde{K}$, it is sufficient to show the following fact: for all $\epsilon > 0$,

$$(5.9) \quad \tilde{K}_{\mathcal{C}}^\top B(\mathcal{C}) \left(B^\top(\mathcal{C}) \tilde{K}_{\mathcal{C}\mathcal{C}} B(\mathcal{C}) + \epsilon \mathbb{I} \right)^{-1} B^\top(\mathcal{C}) \tilde{K}_{\mathcal{C}} \preceq \tilde{K},$$

since by taking the limit $\epsilon \rightarrow 0$, we obtain $\widetilde{L}(\mathcal{C}) \preceq \tilde{K}$. To prove the inequality (5.9), we define $\tilde{K} = AA^\top$ and, thanks to the push-through identity (see Lemma 6.3 in Appendix), we show that

$$A_{\mathcal{C}}^\top B(\mathcal{C}) \left(B^\top(\mathcal{C}) A_{\mathcal{C}} A_{\mathcal{C}}^\top B(\mathcal{C}) + \epsilon \mathbb{I} \right)^{-1} B^\top(\mathcal{C}) A_{\mathcal{C}} = \left(A_{\mathcal{C}}^\top B(\mathcal{C}) B^\top(\mathcal{C}) A_{\mathcal{C}} + \epsilon \mathbb{I} \right)^{-1} A_{\mathcal{C}}^\top B(\mathcal{C}) B^\top(\mathcal{C}) A_{\mathcal{C}} \preceq \mathbb{I},$$

⁵The projected Nyström approximation is denoted by a \tilde{L} (for Low rank matrix), although it is not used in the construction of an extended L -ensemble. The difference should be clear from the context.

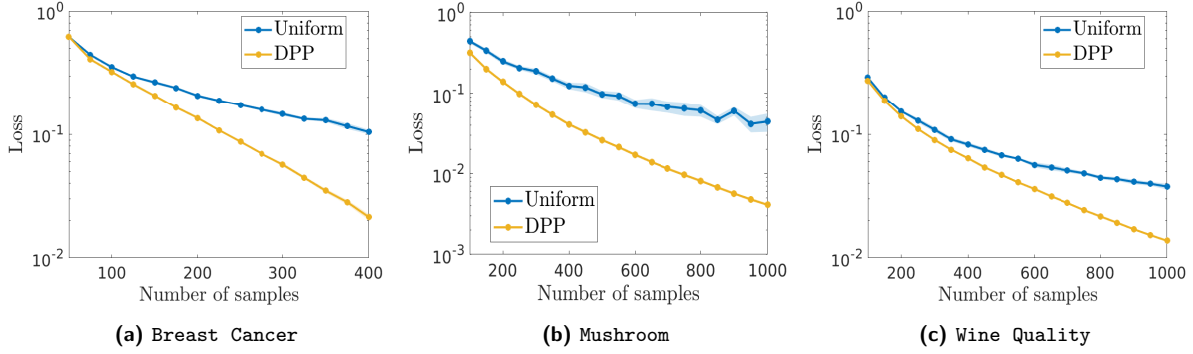


Fig. 4: The relative Nyström approximation error (5.10) as a function of the number landmarks using uniform vs extended L -ensemble DPP sampling.

where $A_C = C^\top A$. (ii) The second identity follows from Proposition 5.1. Consider first the case $\lambda = 1$ without loss of generality. The expectation of the projected Nyström approximation reads

$$\mathbb{E}[\widetilde{L}(\mathcal{C})] = \widetilde{K} \mathbb{E}_C \left[CB(\mathcal{C}) \left(B^\top(\mathcal{C}) \widetilde{K}_{CC} B(\mathcal{C}) \right)^{-1} B^\top(\mathcal{C}) C^\top \right] \widetilde{K} = \widetilde{K} (\widetilde{K} + \mathbb{I})^{-1} \widetilde{K},$$

where Proposition 5.1 was used for the last equality. This gives $\widetilde{K} - \mathbb{E}[\widetilde{L}(\mathcal{C})] = \widetilde{K} (\widetilde{K} + \mathbb{I})^{-1}$. The final result is obtained by replacing \widetilde{K} by \widetilde{K}/λ . \square

Corollary 5.6 only considers matrices which are only non trivial in the subspace V_\perp . Interestingly, the projected Nyström approximation can be written as $\widetilde{L}(\mathcal{C}) = \mathbb{P}_{V_\perp} L(\mathcal{C}) \mathbb{P}_{V_\perp}$, where, in the notations of Proposition 5.1, the ‘unprojected’ Nyström approximation writes

$$L(\mathcal{C}) = KI(\mathcal{C})K.$$

The difference between K and the latter matrix also satisfies simple identities given below.

COROLLARY 5.7. *Let $\mathcal{C} \sim \text{DPP}(K/\lambda, V)$ with $\lambda > 0$. Then, we have*

$$\begin{aligned} \mathbb{E}_C[K - L(\mathcal{C})] &= K - K(\widetilde{K} + \lambda \mathbb{P}_{V_\perp})^+ K \\ \mathbb{E}_C[(CV_C)^+(K - L(\mathcal{C}))] &= V^+ \mathbb{E}_C[K - L(\mathcal{C})] \\ \mathbb{E}_C[(CV_C)^+(K - L(\mathcal{C}))(CV_C)^{\top}] &= \lambda V^+ V^{\top} + V^+ \mathbb{E}_C[K - L(\mathcal{C})] V^{\top}. \end{aligned}$$

Additional properties of the projected Nyström approximation error $\widetilde{K} - \widetilde{L}(\mathcal{C})$ can be obtained by merely replacing K by \widetilde{K} in the above expressions.

Proof. The results simply follow from the identities obtained respectively by: (i) multiplying (5.1) by K on the left and on the right, (ii) multiplying (5.2) by K on the right, (iii) reformulating (5.3) by using (5.1). \square

Empirical results for matrix Nyström approximation. We illustrate here the effect of extended L -ensemble sampling versus uniform sampling for Nyström matrix approximation on the UCI benchmark data sets⁶: **Breast Cancer**, **Mushroom** and **Wine Quality**. The data sets are standardized, and the performance is measured by the relative Frobenius norm of the error:

$$(5.10) \quad \|\widetilde{K} - \widetilde{L}(\mathcal{C})\|_F / \|\widetilde{K}\|_F.$$

We again use a Gaussian kernel $k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|_2^2 / \sigma^2)$ and linear regression component: $V = [X \mathbf{1}_n]$ where $X = [\mathbf{x}_1 \dots \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$. The bandwidth is determined using the median heuristic [34]

⁶<https://archive.ics.uci.edu/ml/index.php>

defined in Section 4.2. The simulation is repeated 10 times and the error bars show the 97.5% confidence interval. The results are displayed in Figure 4. Extended DPP L -ensemble sampling gives a more accurate Nyström approximation. We emphasize that these results are illustrative and that we do not claim that the aforementioned semi-parametric setting is the most suitable for these three data sets. The Nyström approximation of the penalized regression problem can now be discussed by using the matrix Nyström approximation that we just introduced.

5.4. Nyström approximation of regularized regression. Given a subset $\mathcal{C} = \{i_1, \dots, i_k\} \subset [n]$, the Nyström approximation allows to reduce the number of parameters from $n+p$ to $k+p$ without overlooking data points. To do so in the setting of this paper, we propose to solve a simplified problem which differs from (PLS) by the domain of the minimization, i.e., we introduce the following problem

$$(\text{NysPLS}) \quad \min_{f \in \mathcal{H}_N} \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 + \gamma J(f),$$

where the domain is defined by

$$\mathcal{H}_N \triangleq \left\{ f(\mathbf{x}) = \sum_{\ell=1}^k \alpha_{i_\ell} k(\mathbf{x}, \mathbf{x}_{i_\ell}) + \sum_{m=1}^p \beta_m p_m(\mathbf{x}) \text{ s.t. } \sum_{\ell=1}^k \alpha_{i_\ell} p_m(\mathbf{x}_{i_\ell}) = 0 \text{ for all } 1 \leq m \leq p \right\}$$

with $k = |\mathcal{C}|$. The domain of the optimization problem (NysPLS) now includes only finite linear combinations of $k(\cdot, \mathbf{x}_{i_\ell})$ for $i_\ell \in \mathcal{C}$, with a specific condition on the coefficients, whereas the domain of the ‘full’ optimization problem (PLS) includes possibly infinite linear combinations. In analogy with (3.2), this condition yields afterwards the constraint $V_{\mathcal{C}}^\top \boldsymbol{\alpha}' = 0$ where $\boldsymbol{\alpha}' = [\alpha_{i_1} \dots \alpha_{i_k}]^\top \in \mathbb{R}^k$. Here, $V = [p_m(\mathbf{x}_i)]_{im}$ is a $n \times p$ matrix.

The solution of (NysPLS) involves a $(k-p) \times (k-p)$ linear system that we write below, after introducing useful notations. Let $B(\mathcal{C}) \in \mathbb{R}^{k \times (k-p)}$ be a matrix whose columns are an orthonormal basis of $(V_{\mathcal{C}})^\perp$, and that is such that $\mathbb{P}_{V_{\mathcal{C}}^\perp} = B(\mathcal{C})B^\top(\mathcal{C})$. Then, after elementary manipulations, the system yields

$$(5.11) \quad \begin{aligned} \boldsymbol{\alpha}^* &= B(\mathcal{C}) \left(B^\top(\mathcal{C}) K_{\mathcal{C}} \mathbb{P}_{V^\perp} K_{\mathcal{C}}^\top B(\mathcal{C}) + n\gamma B^\top(\mathcal{C}) K_{\mathcal{C}\mathcal{C}} B(\mathcal{C}) \right)^{-1} B^\top(\mathcal{C}) K_{\mathcal{C}} \mathbb{P}_{V^\perp} \mathbf{y}, \\ \boldsymbol{\beta}^* &= (V^\top V)^{-1} V^\top (\mathbf{y} - K_{\mathcal{C}}^\top B(\mathcal{C}) B^\top(\mathcal{C}) \boldsymbol{\alpha}^*). \end{aligned}$$

The details of this derivation are given in Appendix. The above linear system involves positive definite matrices and can be solved conveniently by the conjugate gradient algorithm, possibly after a preconditioning which is described in Appendix.

In-sample estimator. Also, it is straightforward to determine the in-sample estimator $\hat{\mathbf{z}}_N = K_{\mathcal{C}}^\top \boldsymbol{\alpha}^* + V^\top \boldsymbol{\beta}^*$, which is then given, in terms of the projected Nyström approximation, as follows

$$\hat{\mathbf{z}}_N = \widetilde{L}(\mathcal{C}) \left(\widetilde{L}(\mathcal{C}) + n\gamma \mathbb{I}_n \right)^{-1} \mathbb{P}_{V^\perp} \mathbf{y} + \mathbb{P}_V \mathbf{y}.$$

Importantly, this estimator can be formally obtained from the estimator of the ‘full’ problem by replacing \widetilde{K} by $\widetilde{L}(\mathcal{C})$. Notice that the computation of $\boldsymbol{\beta}^*$ only requires solving a $p \times p$ linear system.

5.5. Bound on the expected risk. Recall our data assumption $y_i = f(\mathbf{x}_i) + \epsilon_i$ where ϵ_i denotes i.i.d. $\mathcal{N}(0, \sigma^2)$ noise with $1 \leq i \leq n$. For convenience, define $z_i = f(\mathbf{x}_i)$ for all $1 \leq i \leq n$. The expected risk of the full regression problem is $\mathcal{R}(\hat{\mathbf{z}}) = \mathbb{E}_\epsilon \|\hat{\mathbf{z}} - \mathbf{z}\|_2^2$. We find an upper bound for the expected risk obtained thanks to the projected Nyström approximation, i.e., $\mathcal{R}(\hat{\mathbf{z}}_N) = \mathbb{E}_\epsilon \|\hat{\mathbf{z}}_N - \mathbf{z}\|_2^2$. In the spirit of the kernel ridge regression and Theorem 2.5 in [30], we can prove the following stability bound on expectation.

THEOREM 5.8 (stability of the expected risk). *Let $\mathcal{C} \sim \text{DPP}(K/\lambda, V)$ with $\lambda > 0$. Then, it holds that*

$$\mathbb{E}_{\mathcal{C}} \left[\sqrt{\frac{\mathcal{R}(\hat{\mathbf{z}}_N)}{\mathcal{R}(\hat{\mathbf{z}})}} \right] \leq 1 + \frac{\lambda}{n\gamma} d_{\text{eff}}(\widetilde{K}/\lambda), \text{ with } d_{\text{eff}}(\widetilde{K}/\lambda) = \text{Tr} \left(\widetilde{K} (\widetilde{K} + \lambda \mathbb{I}_n)^{-1} \right).$$

Proof. The proof follows exactly the same lines as in [30, Theorem 3], where \widetilde{K} and $\widetilde{L}(\mathcal{C})$ replace the *psd* kernel matrix and its common Nyström approximation. This can be done since $\widetilde{L}(\mathcal{C})$ formally satisfies the same identity as the common Nyström approximation in kernel ridge regression case when the sampling is done with a L -ensemble. \square

This result indicates that using the Nyström approximation cannot dramatically deteriorate the risk on expectation. An analogous result for leverage scores sampling holding with high probability can be found in [46] for kernel ridge regression. We remark that the effective dimension $d_{\text{eff}}(\tilde{K}/\lambda)$ is crucial in many sampling methods (see also, e.g. [26]). Typically, a small $\lambda > 0$ yields a large expected sample $\mathbb{E}[|C|] = p + d_{\text{eff}}(\tilde{K}/\lambda)$ and therefore reduces the magnitude of the upper bound in Theorem 5.8.

5.6. Application: non-linear time series using semi-parametric models. A typical (embedded) application that requires a small number of parameters is non-linear time series estimation which is a problem of interest in engineering, for instance in the context of system identification, electromechanical systems or for the control of chemical processes. Within the framework of non-linear time series, a common approach consists in estimating a non-linear black-box model to produce accurate and fast forecasts starting from a set of observations. The user usually has some expert knowledge to incorporate into the estimation. This makes the use of semi-parametric regression models especially appealing for systems and control [27, 28], while we refer to [31] for an overview of the various applications in finance, climate and environment sciences. Empirically, it is common to transform time series estimation or system identification problems into regression problems, such as (PLS), as we explain below. Therefore, we use this engineering application as a case study for the function estimation framework used in this paper. A recent theoretical analysis of this type of regression frameworks for time series can be found in [41].

The time cruciality of industrial applications necessitates models with a small number of parameters, as these have a large impact on the memory requirements and prediction speed in real-time forecasting. We demonstrate the use of DPP sampling for Nyström-based regression. We show that Nyström-based regression (NysPLS) does not have a much lower performance compared to the full system, with a lower memory cost and prediction time.

In this simulation, we compare the performance of solving (PLS) and (NysPLS) by using either uniform or DPP sampling. Each model contains a linear parametric part corresponding to the system to be estimated and non-parametric part based on the Gaussian kernel. The non-parametric part can be viewed as a misspecification error. A simple observation given in Proposition 5.9 justifies a decomposition with a separation of variables between the linear and non-linear components.

PROPOSITION 5.9. *Let $\mathbf{v} \in \mathbb{R}^d$ and let $\mathbb{P}_{\mathbf{v}^\perp}$ be the projector onto the orthogonal of \mathbf{v} . Then, the kernel $k(\mathbb{P}_{\mathbf{v}^\perp}\mathbf{x}, \mathbb{P}_{\mathbf{v}^\perp}\mathbf{x}') = \exp(-\|\mathbb{P}_{\mathbf{v}^\perp}(\mathbf{x} - \mathbf{x}')\|^2)$ is positive semi-definite on \mathbb{R}^d .*

Proof. This can be shown thanks to the following result of [7, Thm 2.2]: $\exp(-g(\mathbf{x}, \mathbf{x}'))$ is positive semi-definite if and only if $g(\mathbf{x}, \mathbf{x}')$ is negative semi-definite with respect to 1. Clearly, $g(\mathbf{x}, \mathbf{x}') = \|\mathbb{P}_{\mathbf{v}^\perp}(\mathbf{x} - \mathbf{x}')\|^2$ satisfies $\sum_{i,j=1}^m \alpha_i \alpha_j g(\mathbf{x}_i, \mathbf{x}_j) \leq 0$ for all finite set of \mathbf{x}_i for $1 \leq i \leq m$ and $\boldsymbol{\alpha} \in \mathbb{R}^m$ such that $\sum_{i=1}^m \alpha_i = 0$. \square

Non-linear time series. Then, in what follows, three systems are defined.

System 1: The first model is a static toy example that is given at time step t with $1 \leq t \leq n$ by

$$y^t = a_2 z_1^t + a_1 + \text{sinc}(x_1^t + x_2^t) + \epsilon^t,$$

with $a_2 = 0.2$, $a_1 = 0.4$ and where superscript t indicates a value obtained at time t . The real y^t is the output of the system at time step t , which is given by the combination of a linear combination of the real input z^t and a non-linear function of two other real inputs x_1^t and x_2^t at time t . The training set is obtained by considering a set of input-output pairs obtained for a sequence of integer time steps $1 \leq t \leq n$. The inputs are sampled independently as follows: x_1 and x_2 are $\mathcal{N}(0, 2)$ random variables, whereas $z \sim \mathcal{N}(0, 2.5)$, and $\epsilon \sim \mathcal{N}(0, 0.05)$ is the noise. The training data for the penalized regression problem are (\mathbf{x}_i, y_i) with⁷

$$\mathbf{x}_i = [z^i \ x_1^i \ x_2^i]^\top \in \mathbb{R}^3 \text{ and } y_i = y^i \text{ for } 1 \leq i \leq n.$$

Remark that the time information is not considered here in order to transform the system identification problem into a regression problem, while non-static systems are given below. Let $\mathbf{x} = [z \ x_1 \ x_2]^\top$. The estimated function is of the form

$$f(\mathbf{x}) = \beta_1 p_1(\mathbf{x}) + \beta_2 p_2(\mathbf{x}) + \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i),$$

⁷Notice that a different font is used for the inputs and output, compared to the (\mathbf{x}_i, y_i) pairs.

where $p_1(\mathbf{x}) = 1$ and $p_2(\mathbf{x}) = x_1$, while the kernel is $k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{\sigma^2} \left((x_2 - x'_2)^2 + (x_3 - x'_3)^2\right)\right)$, where x_k denotes the k -th component of \mathbf{x} with $1 \leq k \leq 3$. The estimation problem is then cast into the form of (PLS) since Proposition 5.9 indicates that $k(\mathbf{x}, \mathbf{x}')$ is *psd*.

System 2: The second model is not static: for integer time steps $1 \leq t \leq n$, it reads

$$y^t = a_1 + a_2 y^{t-1} + a_3 y^{t-2} + 2 \operatorname{sinc}(x_1^t + x_2^t) + \epsilon^t,$$

where $a_1 = 0.3$, $a_2 = 0.2$, $a_3 = 0.1$. It is common to define $y^0 = y^{-1} = 0$. Also, we consider independent random variables $x_1 \sim \mathcal{N}(0, 2)$, $x_2 \sim \mathcal{N}(0, 2)$, and a noise $\epsilon \sim \mathcal{N}(0, 0.05)$. The training data for the penalized regression problem are (\mathbf{x}_i, y_i) for $1 \leq i \leq n$ with

$$\mathbf{x}_i = [x_1^i \ x_2^i \ y^{i-1} \ y^{i-2}]^\top \in \mathbb{R}^4 \text{ and } y_i = y^i.$$

Here, the time series is encoded into the regression problem in such a way that \mathbf{x}_i contains the previous two time steps. Let $\mathbf{x} \in \mathbb{R}^4$. The estimated function is of the form

$$f(\mathbf{x}) = \beta_1 p_1(\mathbf{x}) + \beta_2 p_2(\mathbf{x}) + \beta_3 p_3(\mathbf{x}) + \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i),$$

where $p_1(\mathbf{x}) = 1$, $p_2(\mathbf{x}) = x_3$, $p_3(\mathbf{x}) = x_4$ and $k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{\sigma^2} \left((x_1 - x'_1)^2 + (x_2 - x'_2)^2\right)\right)$.

System 3: The third model is of the form:

$$y^t = a_1 + a_2 y^{t-1} + a_3 y^{t-2} + b_1 \operatorname{sinc}(u^{t-1}) + b_2 \operatorname{sinc}(u^{t-2}) + \epsilon^t,$$

with $a_1 = 0.6$, $a_2 = 0.4$, $a_3 = 0.2$, $b_1 = 0.7$, $b_2 = 0.6$, $u \sim \mathcal{N}(0, 4)$, and the noise $\epsilon \sim \mathcal{N}(0, 0.05)$. We also have $y^0 = y^{-1} = 0$ and $u^0 = u^{-1} = 0$ by definition. The training data for the penalized regression problem are (\mathbf{x}_i, y_i) for $1 \leq i \leq n$ with

$$\mathbf{x}_i = [u^{i-1} \ u^{i-2} \ y^{i-1} \ y^{i-2}]^\top \in \mathbb{R}^4 \text{ and } y_i = y^i.$$

Let $\mathbf{x} \in \mathbb{R}^4$. The estimated function is of the form

$$f(\mathbf{x}) = \beta_1 p_1(\mathbf{x}) + \beta_2 p_2(\mathbf{x}) + \beta_3 p_3(\mathbf{x}) + \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i),$$

where $p_1(\mathbf{x}) = 1$, $p_2(\mathbf{x}) = x_3$, $p_3(\mathbf{x}) = x_4$ and $k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{\sigma^2} \left((x_1 - x'_1)^2 + (x_2 - x'_2)^2\right)\right)$.

Simulation setting. We take $n = 1000$ time steps. The data set is split into a 50/25/25 train, validation and test set. The validation set is used to determine the regularization parameter γ and bandwidth σ . For each model, we measure the parameter identification error, i.e., the mean squared error between the true coefficients a_1, a_2, a_3 of the parametric component and their estimates $\beta_1, \beta_2, \beta_3$. More importantly, we calculate the prediction error on the test set: $(1/n_{\text{test}}) \sum_{i=1}^{n_{\text{test}}} (y_i - \hat{f}(\mathbf{x}_i))^2$. The simulation is repeated 10 times. The results are visualized in Figures 5 and 6 where the error bars show the 97.5% confidence interval. Both sampling algorithms are capable of correctly identifying the linear part of the model. Given a number of landmark points, DPP sampling shows better performance than uniform sampling for the prediction error which is the task of practical interest.

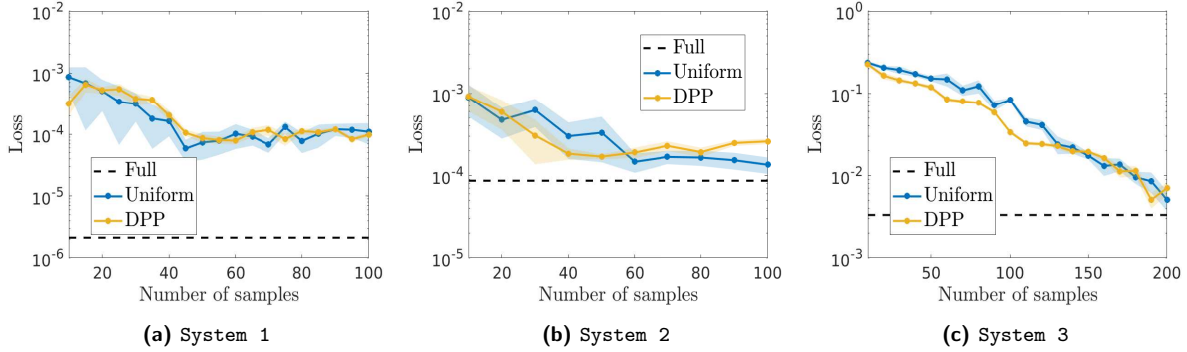


Fig. 5: The parameter identification error as a function of the number landmarks $|\mathcal{C}|$ using uniform vs extended L -ensemble sampling. Here, the total number of training points is $n = 500$.

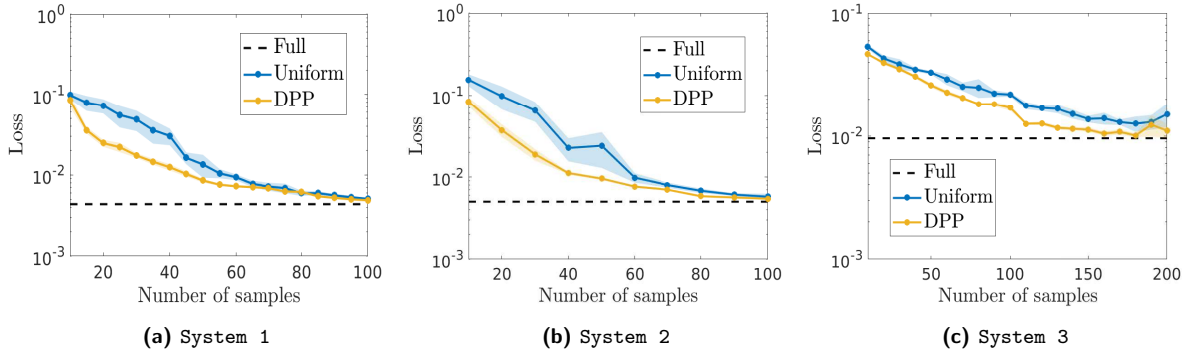


Fig. 6: The prediction error as a function of the number landmarks $|\mathcal{C}|$ using uniform vs extended L -ensemble sampling. The total number of training points is $n = 500$. The loss is the MSE on the test set.

6. Discussion and conclusion. Although the extended L -ensembles are particularly suited for sampling in the context of semi-parametric regression, the sampling cost in practice might be high if an exact DPP sampling algorithm is used. We acknowledge this difficulty and point the reader towards the recent advances in DPP sampling algorithm which have achieved an improved scalability, especially, if the number of sampled landmarks is not large [12]. Sampling a fixed-size DPP without looking at all items was also studied in [9], which provides theoretical guarantees. We expect the methods of [12, 9] to be also applicable for partial-projection DPPs, while further approximate DPP sampling algorithms can be developed in the future.

Acknowledgments. We are grateful to Michał Dereziński for his insightful correspondence about DPPs and implicit regularization. We thank Simon Barthelmé for pointing out the references [42, 38].

Appendix.

6.1. Solution of the penalized least-squares regression.

THEOREM 6.1 (existence, Thm 2.9 in [35]). *Suppose $L(f)$ is a continuous and convex functional in a Hilbert space H and $J(f)$ is a square (semi) norm in H with a null space \mathcal{N}_J , of finite dimension. If $L(f)$ has a unique minimizer in \mathcal{N}_J , then $L(f) + \gamma J(f)$ has a minimizer in H .*

6.2. Useful results.

6.2.1. Classical identities. First, we mention an instrumental result given in [37]. Next, we list a few well-known lemmata.

LEMMA 6.2 (theorem 2.1 in [37]). *Let $\mathcal{A} \subseteq [n]$ and $M \in \mathbb{R}^{n \times n}$. Let $\mathbb{1}_{\bar{\mathcal{A}}}$ be the diagonal matrix with ones in the diagonal positions corresponding to elements of $\bar{\mathcal{A}} = [n] \setminus \mathcal{A}$, and zeros otherwise. Then, it holds that $\sum_{\mathcal{C}: \mathcal{A} \subseteq \mathcal{C}} \det M_{\mathcal{C}\mathcal{C}} = \det(M + \mathbb{1}_{\bar{\mathcal{A}}})$.*

LEMMA 6.3 (push-through). *Let $X \in \mathbb{R}^{m \times k}$ and $Y \in \mathbb{R}^{k \times m}$. Then, it holds that*

$$(XY + \mathbb{I}_m)^{-1}X = X(YX + \mathbb{I}_k)^{-1}.$$

LEMMA 6.4 (Cauchy-Binet). *Let $X \in \mathbb{R}^{n \times p}$ and $Y \in \mathbb{R}^{n \times p}$ with $n > p$. Then, it holds that*

$$\det(X^\top Y) = \sum_{\mathcal{C} \subseteq [n]: |\mathcal{C}|=p} \det(X_{\mathcal{C}}^\top Y_{\mathcal{C}}),$$

where $X_{\mathcal{C}}$ is the $p \times p$ matrix obtained from X by selecting the rows indexed by \mathcal{C} .

LEMMA 6.5. *Let $A \in \mathbb{R}^n$ and $W \in \mathbb{R}^{n \times p}$ such that A is invertible as well as $W^\top A^{-1}W$. Then, it holds that*

$$\begin{pmatrix} A & W \\ W^\top & 0 \end{pmatrix}^{-1} = \begin{pmatrix} A^{-1} - A^{-1}W(W^\top A^{-1}W)^{-1}W^\top A^{-1} & A^{-1}W(W^\top A^{-1}W)^{-1} \\ (W^\top A^{-1}W)^{-1}W^\top A^{-1} & -(W^\top A^{-1}W)^{-1} \end{pmatrix}.$$

We also use a slightly different result.

LEMMA 6.6. *Let $A \in \mathbb{R}^{n \times n}$. Let $W \in \mathbb{R}^{n \times p}$ be a matrix with full column rank. Let B be a matrix with orthonormal columns so that $\mathbb{P}_{W^\perp} = BB^\top$. If $B^\top AB$ is invertible, we have*

$$\begin{pmatrix} A & W \\ W^\top & 0 \end{pmatrix}^{-1} = \begin{pmatrix} (\mathbb{P}_{W^\perp} A \mathbb{P}_{W^\perp})^+ & (\mathbb{I} - (\mathbb{P}_{W^\perp} A \mathbb{P}_{W^\perp})^+) W^{+\top} \\ W^+ (\mathbb{I} - A(\mathbb{P}_{W^\perp} A \mathbb{P}_{W^\perp})^+) & -W^+ (A - A(\mathbb{P}_{W^\perp} A \mathbb{P}_{W^\perp})^+) W^{+\top} \end{pmatrix},$$

where $W^+ = (W^\top W)^{-1}W^\top$ and $(\mathbb{P}_{W^\perp} A \mathbb{P}_{W^\perp})^+ = B(B^\top AB)^{-1}B^\top$.

LEMMA 6.7. *Let M be a $k \times k$ matrix and let S be a $n \times k$ matrix such that $S^\top S = \mathbb{I}_k$. Then, we have $(SMS^\top)^+ = SM^+S^\top$.*

Proof. The criteria satisfied by the Moore-Penrose pseudo-inverse are readily checked. It holds that

$$\begin{aligned} (SMS^\top)(SMS^\top)^+(SMS^\top) &= SM^+MM^+S^\top = SMS^\top \\ (SMS^\top)^+(SMS^\top)(SMS^\top)^+ &= SMM^+MS^\top = SM^+S^\top, \end{aligned}$$

while the following matrices are Hermitian

$$\begin{aligned} (SMS^\top)(SMS^\top)^+ &= SM^+MS^\top \\ (SMS^\top)^+(SMS^\top) &= SMM^+S^\top, \end{aligned}$$

since M^+M and MM^+ are Hermitian. □

LEMMA 6.8 (matrix determinant lemma). *Let $A \in \mathbb{R}^{n \times n}$ be an invertible matrix and $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$. Then, it holds that $\det(A + \mathbf{u}\mathbf{v}^\top) = (1 + \mathbf{v}^\top A^{-1}\mathbf{u}) \det(A)$.*

6.2.2. Lemmata related to the extended L -ensemble formalism. We below provide helpful formulae to calculate the determinant of matrices with the special structure of extended L -ensembles.

LEMMA 6.9 (Lemma 3.12 in [2]). *Let (M, V) is a NNP such as in Definition 2.1. Let $\widetilde{M} = \mathbb{P}_{V^\perp} M \mathbb{P}_{V^\perp} \in \mathbb{R}^{n \times n}$ and $V = QR \in \mathbb{R}^{n \times p}$ where $Q \in \mathbb{R}^{n \times p}$ has orthonormal columns and R is upper triangular. Then, we have*

$$\det \begin{pmatrix} \widetilde{M} & Q \\ Q^\top & 0 \end{pmatrix} = (-1)^p [t^p] \det \left(\widetilde{M} + tQQ^\top \right),$$

where $[t^p]$ denotes the coefficient of the term t^p .

LEMMA 6.10 (lemma 3.11 [2]). *Let $M \in \mathbb{R}^{n \times n}$ and $V \in \mathbb{R}^{n \times p}$ such that (M, V) is a NNP. Then, we have*

$$\det \begin{pmatrix} M & V \\ V^\top & 0 \end{pmatrix} = (-1)^p \det(V^\top V) \det(Q_\perp^\top M Q_\perp),$$

where $Q_\perp \in \mathbb{R}^{n \times (n-p)}$ has orthonormal columns and is such that $V^\top Q_\perp = 0$.

For the sake of completeness, we derive an equivalent expression for the normalization appearing in Definition 2.3, by extensively using proof techniques of [2, 4].

PROPOSITION 6.11 (normalization factors). *Let (L, V) is a NNP such as in Definition 2.1. We have the following identity*

$$(-1)^p \det(\mathbb{I} + \tilde{L}) \det(V^\top V) = \det \begin{pmatrix} L + \mathbb{I} & V \\ V^\top & 0 \end{pmatrix}.$$

Proof. Let $V = QR$ thanks to the QR-decomposition of V , with $Q \in \mathbb{R}^{n \times p}$. Let $Q_\perp \in \mathbb{R}^{n \times (n-p)}$ be a matrix with orthonormal columns so that $\mathbb{P}_{V^\perp} = Q_\perp Q_\perp^\top$. Then, $\tilde{L} = Q_\perp Q_\perp^\top L Q_\perp Q_\perp^\top$. This gives the following decomposition

$$(6.1) \quad \mathbb{I} + \tilde{L} = QQ^\top + Q_\perp (\mathbb{I}_{n-p} + Q_\perp^\top L Q_\perp) Q_\perp^\top,$$

where we used $\mathbb{I} = Q_\perp Q_\perp^\top + QQ^\top$. Define the eigendecomposition $Q_\perp^\top L Q_\perp + \mathbb{I}_{n-p} = U \Lambda U^\top$, with $U \in \mathbb{R}^{(n-p) \times (n-p)}$ an orthogonal matrix and $\Lambda \in \mathbb{R}^{(n-p) \times (n-p)}$ a diagonal matrix. Therefore, in the light of (6.1), we find $\det(\mathbb{I} + \tilde{L}) = \det(\Lambda)$. Now, we use the following identity, which results from Lemma 2.6 in [4],

$$(6.2) \quad \det \begin{pmatrix} L + \mathbb{I} & V \\ V^\top & 0 \end{pmatrix} = \det \begin{pmatrix} \tilde{L} + \tilde{\mathbb{I}} & V \\ V^\top & 0 \end{pmatrix},$$

with $\tilde{\mathbb{I}} = \mathbb{P}_{V^\perp}$ and $\tilde{L} = \mathbb{P}_{V^\perp} L \mathbb{P}_{V^\perp}$. Next, by using the decomposition of $V = QR$, we find

$$(6.3) \quad \det \begin{pmatrix} \tilde{L} + \tilde{\mathbb{I}} & V \\ V^\top & 0 \end{pmatrix} = \det(R^\top R) \det \begin{pmatrix} \tilde{L} + \tilde{\mathbb{I}} & Q \\ Q^\top & 0 \end{pmatrix} = \det(V^\top V) \det \begin{pmatrix} \tilde{L} + \tilde{\mathbb{I}} & Q \\ Q^\top & 0 \end{pmatrix},$$

where the first equality uses Lemma 6.10. Then, we use Lemma 6.9 to express the last factor of the RHS of (6.3),

$$(6.4) \quad \det \begin{pmatrix} \tilde{L} + \tilde{\mathbb{I}} & Q \\ Q^\top & 0 \end{pmatrix} = (-1)^p [t^p] \det \left(\tilde{L} + \tilde{\mathbb{I}} + tQQ^\top \right),$$

where $[t^p]p(t)$ denotes the coefficient of the term t^p in the polynomial $p(t)$. Next, by using the above eigendecomposition, we find

$$(6.5) \quad [t^p] \det \left(\tilde{L} + \tilde{\mathbb{I}} + tQQ^\top \right) = [t^p] \det \left[\begin{pmatrix} Q & U \end{pmatrix} \begin{pmatrix} t\mathbb{I}_{p \times p} & 0 \\ 0 & \Lambda \end{pmatrix} \begin{pmatrix} Q^\top \\ U^\top \end{pmatrix} \right] = \det(\Lambda) = \det(\mathbb{I} + \tilde{L}),$$

where we used that $\begin{pmatrix} Q & U \end{pmatrix} \in \mathbb{R}^{n \times n}$ is an orthogonal matrix. Finally, we combine (6.2), (6.3), (6.4) and (6.5) to give

$$\det \begin{pmatrix} L + \mathbb{I} & V \\ V^\top & 0 \end{pmatrix} = \det(V^\top V) (-1)^p \det(\mathbb{I} + \tilde{L}),$$

which is the desired result. \square

Below, we give a lemma allowing to simplify several expressions in the paper.

LEMMA 6.12 ([4]). *Let (K, V) be a NNP. Let $B(\mathcal{C}) \in \mathbb{R}^{k \times (k-p)}$ be a matrix whose columns are an orthonormal basis of $(V_{\mathcal{C}})^\perp$. Let $\tilde{K} = \mathbb{P}_{V^\perp} K \mathbb{P}_{V^\perp}$. Then, it holds that $\mathbb{P}_V C B(\mathcal{C}) = 0$ and*

$$B^\top(\mathcal{C}) K_{\mathcal{C}\mathcal{C}} B(\mathcal{C}) = B^\top(\mathcal{C}) \tilde{K}_{\mathcal{C}\mathcal{C}} B(\mathcal{C}).$$

Proof. By definition, $\mathbb{P}_{V^\perp} = \mathbb{I} - QQ^\top$ where Q has orthonormal columns and is obtained thanks to the QR-decomposition $V = QR$. Then, we find $B^\top(\mathcal{C}) V_{\mathcal{C}} = 0 = B^\top(\mathcal{C}) Q_{\mathcal{C}}$. Therefore, we obtain $B^\top(\mathcal{C}) \tilde{K}_{\mathcal{C}\mathcal{C}} B(\mathcal{C}) = B^\top(\mathcal{C}) K_{\mathcal{C}\mathcal{C}} B(\mathcal{C})$ by using $\tilde{K}_{\mathcal{C}\mathcal{C}} = (C^\top - Q_{\mathcal{C}} Q_{\mathcal{C}}^\top) K (C - Q Q_{\mathcal{C}}^\top)$, and $B^\top(\mathcal{C}) Q_{\mathcal{C}} = 0$. \square

6.3. Proof of Theorem 4.1 . Let $\mathbf{u}, \mathbf{v} \in \mathbb{R}^{n+p}$. First, we first notice that

$$\begin{pmatrix} C^\top K C & C^\top V \\ V^\top C & 0 \end{pmatrix} = \begin{pmatrix} C^\top & 0 \\ 0 & \mathbb{I} \end{pmatrix} \begin{pmatrix} K & V \\ V^\top & 0 \end{pmatrix} \begin{pmatrix} C & 0 \\ 0 & \mathbb{I} \end{pmatrix} = \tilde{C}^\top \begin{pmatrix} K & V \\ V^\top & 0 \end{pmatrix} \tilde{C},$$

where \tilde{C} is a sampling matrix corresponding to a subset of $\{1, \dots, n+p\}$, i.e., \tilde{C} is associated to the set $\tilde{\mathcal{C}} = \mathcal{C} \cup \mathcal{A}$ with $\mathcal{A} = \{n+1, \dots, n+p\}$. Therefore, by using Lemma 6.2, we have the following identity $\sum_{\tilde{\mathcal{C}}: \mathcal{A} \subseteq \tilde{\mathcal{C}}} \det Q_{\tilde{\mathcal{C}}\tilde{\mathcal{C}}} = \det(Q + \mathbb{1}_{\mathcal{A}})$ for all $Q \in \mathbb{R}^{(n+p) \times (n+p)}$. In particular, for $Q = \begin{pmatrix} K & V \\ V^\top & 0 \end{pmatrix}$, this gives

$$(6.6) \quad \sum_{\mathcal{C} \subseteq [n]} \det \begin{pmatrix} C^\top K C & C^\top V \\ V^\top C & 0 \end{pmatrix} = \det \begin{pmatrix} K + \mathbb{I} & V \\ V^\top & 0 \end{pmatrix}.$$

Then, the desired expectation can be written as follows

$$\frac{1}{N} \sum_{\mathcal{C} \subseteq [n]} \det(\tilde{C}^\top Q \tilde{C}) \times \mathbf{u}^\top \tilde{C} \left(\tilde{C}^\top Q \tilde{C} \right)^{-1} \tilde{C}^\top \mathbf{v} = \frac{1}{N} \sum_{\mathcal{C} \subseteq [n]} \det(\tilde{C}^\top Q \tilde{C}) - \frac{1}{N} \sum_{\mathcal{C} \subseteq [n]} \det(\tilde{C}^\top (Q - \mathbf{v}\mathbf{u}^\top) \tilde{C}),$$

thanks to the matrix determinant lemma (Lemma 6.8), and where the normalization N is given by (6.6). Define for simplicity $T_1 \triangleq \sum_{\mathcal{C} \subseteq [n]} \det(\tilde{C}^\top Q \tilde{C})$ and $T_2 \triangleq \sum_{\mathcal{C} \subseteq [n]} \det(\tilde{C}^\top (Q - \mathbf{v}\mathbf{u}^\top) \tilde{C})$. Then, we find

$$T_1 = \det \begin{pmatrix} K + \mathbb{I} & V \\ V^\top & 0 \end{pmatrix} = \det(Q + \mathbb{1}_{[n]}) = N,$$

where $\mathbb{1}_{\mathcal{B}}$ the diagonal matrix with ones in the diagonal positions corresponding to elements of the set \mathcal{B} , and zeros otherwise. Let $\mathbf{u} = [\mathbf{u}_0, \mathbf{u}_1]^\top$ and $\mathbf{v} = [\mathbf{v}_0, \mathbf{v}_1]^\top$, with $\mathbf{u}_0, \mathbf{v}_0 \in \mathbb{R}^n$ and $\mathbf{u}_1, \mathbf{v}_1 \in \mathbb{R}^p$. Similarly, we have also

$$\begin{aligned} T_2 &= \sum_{\mathcal{C} \subseteq [n]} \det \begin{pmatrix} C^\top (K - \mathbf{v}_0 \mathbf{u}_0^\top) C & C^\top (V - \mathbf{v}_0 \mathbf{u}_1^\top) \\ (V^\top - \mathbf{v}_1 \mathbf{u}_0^\top) C & -\mathbf{v}_1 \mathbf{u}_1^\top \end{pmatrix} = \det \begin{pmatrix} (K + \mathbb{I} - \mathbf{v}_0 \mathbf{u}_0^\top) & V - \mathbf{v}_0 \mathbf{u}_1^\top \\ V^\top - \mathbf{v}_1 \mathbf{u}_0^\top & -\mathbf{v}_1 \mathbf{u}_1^\top \end{pmatrix} \\ &= \det(Q + \mathbb{1}_{[n]} - \mathbf{v}\mathbf{u}^\top). \end{aligned}$$

Hence, we obtain another expression for the desired expectation

$$\begin{aligned} \mathbb{E}_{\mathcal{C} \sim DPP(K, V)} \left[\begin{pmatrix} \mathbf{u}_0, \mathcal{C} \\ \mathbf{u}_1 \end{pmatrix}^\top \begin{pmatrix} C^\top K C & C^\top V \\ V^\top C & 0 \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{v}_0, \mathcal{C} \\ \mathbf{v}_1 \end{pmatrix} \right] &= \frac{\det(Q + \mathbb{1}_{[n]}) - \det(Q + \mathbb{1}_{[n]} - \mathbf{v}\mathbf{u}^\top)}{\det(Q + \mathbb{1}_{[n]})} \\ &= \mathbf{u}^\top (Q + \mathbb{1}_{[n]})^{-1} \mathbf{v}, \end{aligned}$$

where the last equality uses the matrix determinant lemma. This completes the proof.

6.4. Details of the derivation of the large scale system of Section 5.4. After elementary manipulations, we obtain the following system

$$\begin{aligned} \left(\mathbb{P}_{V_c^\perp} K_C \mathbb{P}_{V^\perp} K_C^\top \mathbb{P}_{V_c^\perp} + n\gamma \mathbb{P}_{V_c^\perp} K_{CC} \mathbb{P}_{V_c^\perp} \right) \boldsymbol{\alpha}' &= \mathbb{P}_{V_c^\perp} K_C \mathbb{P}_{V^\perp} \mathbf{y} \\ \boldsymbol{\beta}' &= (V^\top V)^{-1} V^\top \left(\mathbf{y} - K_C^\top \mathbb{P}_{V_c^\perp} \boldsymbol{\alpha}' \right) \\ \mathbb{P}_{V_c} \boldsymbol{\alpha}' &= 0, \end{aligned}$$

which yields the system in (5.11), as we show below. Let $B(\mathcal{C}) \in \mathbb{R}^{k \times (k-p)}$ be a matrix whose columns are an orthonormal basis of $(V_c)^\perp$, and that is such that $\mathbb{P}_{V_c^\perp} = B(\mathcal{C}) B^\top(\mathcal{C})$. Define $\Xi = B^\top(\mathcal{C}) K_{CC} B(\mathcal{C})$, which is non-singular almost surely, as shown in the proof of Proposition 5.1. Then, the first equation of the system yields

$$\begin{aligned} \boldsymbol{\alpha}' &= B(\mathcal{C}) \left(B^\top(\mathcal{C}) K_C \mathbb{P}_{V^\perp} K_C^\top B(\mathcal{C}) + n\gamma B^\top(\mathcal{C}) K_{CC} B(\mathcal{C}) \right)^{-1} B^\top(\mathcal{C}) K_C \mathbb{P}_{V^\perp} \mathbf{y} \\ &= B(\mathcal{C}) \Xi^{-1/2} \left(\Xi^{-1/2} B^\top(\mathcal{C}) K_C \mathbb{P}_{V^\perp} K_C^\top B(\mathcal{C}) \Xi^{-1/2} + n\gamma \mathbb{I}_{k-p} \right)^{-1} \Xi^{-1/2} B^\top(\mathcal{C}) K_C \mathbb{P}_{V^\perp} \mathbf{y} \\ &= B(\mathcal{C}) \Xi^{-1} B^\top(\mathcal{C}) K_C \mathbb{P}_{V^\perp} \left(\mathbb{P}_{V^\perp} K_C^\top B(\mathcal{C}) \Xi^{-1} B^\top(\mathcal{C}) K_C \mathbb{P}_{V^\perp} + n\gamma \mathbb{I}_n \right)^{-1} \mathbb{P}_{V^\perp} \mathbf{y}, \end{aligned}$$

where we used the push-through identity for the last equality (Lemma 6.3) $(XY + \mathbb{I})^{-1}X = X(YX + \mathbb{I})^{-1}$ with $X = \Xi^{-1/2}B^\top(\mathcal{C})K_{\mathcal{C}}\mathbb{P}_{V^\perp}$, $Y = X^\top$ and $\mathbb{P}_{V^\perp}^2 = \mathbb{P}_{V^\perp}$. As detailed above in Lemma 6.12, an equivalent expression for $\Xi = B^\top(\mathcal{C})K_{\mathcal{C}}B(\mathcal{C})$ is $\Xi = B^\top(\mathcal{C})\tilde{K}_{\mathcal{C}}B(\mathcal{C})$. The in-sample estimator $\hat{\mathbf{z}}_N = K_{\mathcal{C}}^\top \boldsymbol{\alpha}'^* + V^\top \boldsymbol{\beta}'^*$ is then given, in terms of the projected Nyström approximation, as follows

$$\hat{\mathbf{z}}_N = \widetilde{L(\mathcal{C})} \left(\widetilde{L(\mathcal{C})} + n\gamma\mathbb{I}_n \right)^{-1} \mathbb{P}_{V^\perp} \mathbf{y} + \mathbb{P}_V \mathbf{y},$$

which is the result stated in Section 5.4.

Preconditioning. Consider the linear system in (5.11) and notice that the largest eigenvalue of the matrix

$$B^\top(\mathcal{C})K_{\mathcal{C}}\mathbb{P}_{V^\perp}K_{\mathcal{C}}^\top B(\mathcal{C}) = B^\top(\mathcal{C})C^\top \tilde{K}^2 CB(\mathcal{C}),$$

in the first term can be possibly numerically large, since it involves \tilde{K}^2 . Therefore, the linear system might be ill-conditioned. A preconditioning may improve the convergence of a linear solver such as the conjugate gradient method. We define the preconditioner as follows. Denote the marginal probabilities of $DPP(K/\lambda, V)$ by

$$(6.7) \quad \boldsymbol{\ell} = \text{diag} \left(\mathbb{P}_V + \tilde{K}(\tilde{K} + \lambda\mathbb{I})^{-1} \right),$$

as given by (1.3), and define the diagonal matrix $D = \text{Diag}(\boldsymbol{\ell})^{-1}$. Then, it simply holds that $\mathbb{E}_{\mathcal{C}}[CD_{\mathcal{C}}C^\top] = \mathbb{I}$, since the marginal probability is $\ell_i = \Pr(i \in \mathcal{Y})$ for $\mathcal{Y} \sim DPP(K/\lambda, L)$. This remark motivates a preconditioning of the above linear system by approximating \tilde{K}^2 by $\tilde{K}CD_{\mathcal{C}}C^\top\tilde{K}$. Let H be a matrix obtained by the following Cholevski decomposition

$$HH^\top = \left(B^\top(\mathcal{C})\tilde{K}_{\mathcal{C}}D_{\mathcal{C}}\tilde{K}_{\mathcal{C}}B(\mathcal{C}) + n\gamma B^\top(\mathcal{C})\tilde{K}_{\mathcal{C}}B(\mathcal{C}) \right)^{-1}.$$

The equivalent resulting linear system

$$(6.8) \quad \boldsymbol{\alpha}'^* = B(\mathcal{C})H \left(H^\top B^\top(\mathcal{C})K_{\mathcal{C}}\mathbb{P}_{V^\perp}K_{\mathcal{C}}^\top B(\mathcal{C})H + n\gamma H^\top B^\top(\mathcal{C})K_{\mathcal{C}}B(\mathcal{C})H \right)^{-1} H^\top B^\top(\mathcal{C})K_{\mathcal{C}}\mathbb{P}_{V^\perp} \mathbf{y},$$

is likely to have a smaller condition number. This type of conditioning was studied in the case of kernel ridge regression in [51].

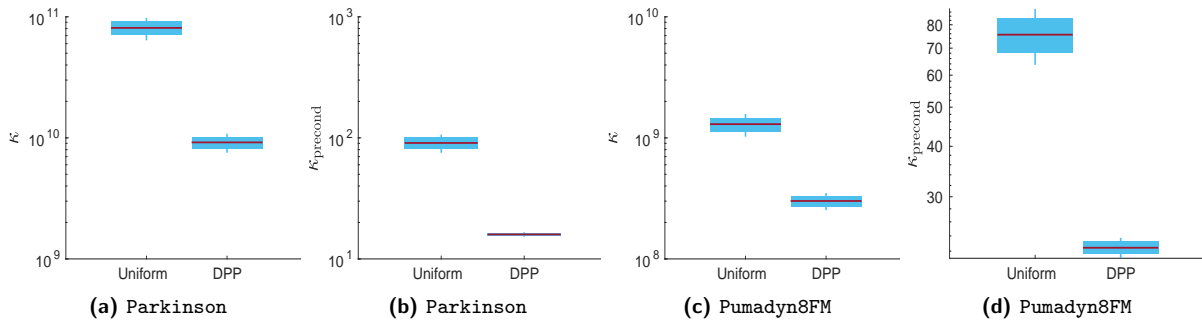


Fig. 7: Preconditioning results. The condition number of the linear system before (κ) and after the preconditioning (κ_{precond}) is plotted for uniform and DPP sampling. From left to right, the condition number before and after preconditioning, for Parkinson and Pumadyn8FM data sets, respectively.

For convenience, we illustrate the use of the preconditioner on the UCI benchmark data sets *Parkinson*, and *Pumadyn8FM*. A Gaussian kernel with $\sigma = 5$ and linear regression component is used after standardizing the data sets: $V = [X \mathbf{1}_n]$ where $X = [\mathbf{x}_1 \dots \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$. We compare the condition number of the

linear system in (5.11) with the preconditioned system using the preconditioner given in (6.8). For the uniform sampling method, we use $D = (n/|\mathcal{C}|)\text{diag}(\mathbf{1}_n)$ in the preconditioner formula as in [51]. The ridge regularization parameter for the linear system as well as the regularization parameter of the DPP are equal to $\lambda = \gamma = 10^{-6}$ for simplicity. The number of samples is equal to the effective dimensionality: $\sum_i \ell_i$, which is 1325 and 636 for the Parkinson, and Pumadyn8FM data sets respectively. The experiment is repeated 10 times. From the results in Figure 7, we empirically see that using the proposed preconditioner in combination with the DPP sampling procedure, results in a smaller condition number of the linear system obtained from (NysPLS).

REFERENCES

- [1] H. AVRON AND C. BOUTSIDIS, *Faster subset selection for matrices and applications*, SIAM Journal on Matrix Analysis and Applications, 34 (2013), pp. 1464–1499.
- [2] S. BARTHELMÉ AND K. USEVICH, *Spectral properties of kernel matrices in the flat limit*, SIAM Journal on Matrix Analysis and Applications, 42 (2021), pp. 17–57.
- [3] S. BARTHELMÉ, P.-O. AMBLARD, AND N. TREMBLAY, *Asymptotic equivalence of fixed-size and varying-size determinantal point processes*, Bernoulli, 25 (2019), pp. 3555–3589.
- [4] S. BARTHELMÉ, N. TREMBLAY, K. USEVICH, AND P.-O. AMBLARD, *Determinantal Point Processes in the Flat Limit: Extended L-ensembles, Partial-Projection DPPs and Universality Classes*, arXiv preprint arXiv:2007.04117, (2020).
- [5] R. BEATSON, W. LIGHT, AND S. BILLINGS, *Fast Solution of the Radial Basis Function Interpolation Equations: Domain Decomposition Methods*, SIAM J. Sci. Comput., 22 (2000), p. 1717–1740.
- [6] A. BELHADJI, R. BARDENET, AND P. CHAINAIS, *A determinantal point process for column subset selection*, Journal of Machine Learning Research, 21 (2020), pp. 1–62.
- [7] C. BERG, J. CHRISTENSEN, AND P. RESSEL, *Harmonic Analysis on Semigroups*, vol. 100 of Graduate Texts in Mathematics, Springer New York, New York, NY, 1984.
- [8] M. BIANCOLINI, *Fast radial basis functions for engineering applications*, Springer, 2017.
- [9] D. CALANDRIELLO, M. DEREZIŃSKI, AND M. VALKO, *Sampling from a k-DPP without looking at all items*, To appear at NeurIPS 2020, preprint arXiv:2006.16947, (2020).
- [10] J. CARR, R. BEATSON, J. CHERRIE, T. MITCHELL, W. FRIGHT, B. MCCALLUM, AND T. EVANS, *Reconstruction and representation of 3D objects with radial basis functions*, in Proceedings of the 28th annual conference on Computer graphics and interactive techniques, 2001, pp. 67–76.
- [11] J. CARR, W. FRIGHT, AND R. BEATSON, *Surface interpolation with radial basis functions for medical imaging*, IEEE transactions on medical imaging, 16 (1997), pp. 96–107.
- [12] M. DEREZIŃSKI, D. CALANDRIELLO, AND M. VALKO, *Exact sampling of determinantal point processes with sublinear time preprocessing*, in Advances in Neural Information Processing Systems, 2019, pp. 11546–11558.
- [13] M. DEREZIŃSKI, K. L. CLARKSON, M. MAHONEY, AND M. WARMUTH, *Minimax experimental design: Bridging the gap between statistical and worst-case approaches to least squares regression*, COLT, (2019).
- [14] M. DEREZIŃSKI, R. KHANNA, AND M. MAHONEY, *Improved guarantees and a multiple-descent curve for the Column Subset Selection Problem and the Nyström method*, To appear at NeurIPS 2020, preprint arXiv:2002.09073, (2020).
- [15] M. DEREZIŃSKI, F. LIANG, Z. LIAO, AND M. MAHONEY, *Precise expressions for random projections: Low-rank approximation and randomized Newton*, To appear at NeurIPS 2020, preprint arXiv:2006.10653, (2020).
- [16] M. DEREZIŃSKI, F. LIANG, AND M. MAHONEY, *Exact expressions for double descent and implicit regularization via surrogate random design*, To appear at NeurIPS 2020, preprint arXiv:1912.04533, (2019).
- [17] M. DEREZIŃSKI, F. LIANG, AND M. MAHONEY, *Bayesian experimental design using regularized determinantal point processes*, in International Conference on Artificial Intelligence and Statistics, 2020, pp. 3197–3207.
- [18] M. DEREZIŃSKI AND M. MAHONEY, *Determinantal Point Processes in Randomized Numerical Linear Algebra*, Notices of the AMS, 68 (2021).
- [19] M. DEREZIŃSKI AND M. WARMUTH, *Unbiased estimates for linear regression via volume sampling*, in Advances in Neural Information Processing Systems, 2017, pp. 3084–3093.
- [20] M. DEREZIŃSKI, M. WARMUTH, AND D. HSU, *Unbiased estimators for random design regression*, preprint arXiv:1907.03411, (2019).
- [21] M. DEREZIŃSKI AND M. K. WARMUTH, *Reverse iterative volume sampling for linear regression*, Journal of Machine Learning Research, 19 (2018), pp. 1–39.
- [22] M. DEREZIŃSKI, M. K. WARMUTH, AND D. HSU, *Correcting the bias in least squares regression with volume-rescaled sampling*, in Proceedings of Machine Learning Research, vol. 89, 2019, pp. 944–953.
- [23] M. DEREZIŃSKI, M. K. WARMUTH, AND D. J. HSU, *Leveraged volume sampling for linear regression*, in Advances in Neural Information Processing Systems, vol. 31, 2018.
- [24] P. DRINEAS, M. MAGDON-ISMAIL, M. MAHONEY, AND D. WOODRUFF, *Fast approximation of matrix coherence and statistical leverage*, Journal of Machine Learning Research, 13 (2012), pp. 3475–3506.
- [25] J. DUCHON, *Splines minimizing rotation-invariant semi-norms in Sobolev spaces*, in Constructive Theory of Functions of Several Variables, 1976.
- [26] A. EL ALAOU AND M. MAHONEY, *Fast randomized kernel ridge regression with statistical guarantees*, in Advances in Neural Information Processing Systems 28, 2015, pp. 775–783.
- [27] M. ESPINOZA, J. SUYKENS, AND B. DE MOOR, *Partially linear models and least squares support vector machines*, in 2004 43rd IEEE Conference on Decision and Control (CDC), vol. 4, 2004, pp. 3388–3393 Vol.4.

- [28] M. ESPINOZA, J. SUYKENS, AND B. DE MOOR, *Kernel based partially linear models and nonlinear identification*, IEEE Transactions on Automatic Control, 50 (2005), pp. 1602–1606.
- [29] M. FANUEL, J. SCHREURS, AND J. SUYKENS, *Nyström landmark sampling and regularized Christoffel functions*, arXiv preprint arXiv:1905.12346, (2019).
- [30] M. FANUEL, J. SCHREURS, AND J. SUYKENS, *Diversity sampling is an implicit regularization for kernel methods*, SIAM Journal on Mathematics of Data Science, 3 (2021), pp. 280–297.
- [31] J. GAO, *Nonlinear time series : semiparametric and nonparametric methods*, vol. 108 of Monographs on statistics and applied probability, Chapman & Hall, New York, 2007.
- [32] B. GAUTHIER, *Approche spectrale pour l'interpolation à noyaux et positivité conditionnelle*, École Nationale Supérieure des Mines de Saint-Étienne, (2011). PhD thesis.
- [33] B. GAUTHIER AND L. PRONZATO, *Spectral approximation of the imse criterion for optimal designs in kernel-based interpolation models*, SIAM/ASA Journal on Uncertainty Quantification, 2 (2014), pp. 805–825.
- [34] A. GRETTON, K. BORGFWARDT, M. RASCH, B. SCHÖLKOPF, AND A. SMOLA, *A kernel two-sample test*, The Journal of Machine Learning Research, 13 (2012), pp. 723–773.
- [35] C. GU, *Smoothing Spline Anova Models*, Springer, 2013.
- [36] R. HUANG AND C. SZEPESVARI, *A Finite-Sample Generalization Bound for Semiparametric Regression: Partially Linear Models*, vol. 33 of Proceedings of Machine Learning Research, 2014, pp. 402–410.
- [37] A. KULESZA AND B. TASKAR, *Determinantal Point Processes for Machine Learning*, Foundations and Trends in Machine Learning, 5 (2012), pp. 123–286.
- [38] R. LYONS, *Determinantal probability measures*, Publ. Math., 98 (2003), p. 167–212.
- [39] M. MAHONEY, *Approximate Computation and Implicit Regularization for Very Large-Scale Data Analysis*, in Proceedings of the 31st ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, 2012, p. 143–154.
- [40] M. MAHONEY AND L. ORECCHIA, *Implementing Regularization Implicitly via Approximate Eigenvector Computation*, in Proceedings of the 28th International Conference on Machine Learning, 2011, p. 121–128.
- [41] Z. MARIET AND V. KUZNETSOV, *Foundations of sequence-to-sequence modeling for time series*, vol. 89 of Proceedings of Machine Learning Research, PMLR, 2019, pp. 408–417.
- [42] S. B. MAURER, *Matrix generalizations of some theorems on trees, cycles and cocycles in graphs*, SIAM Journal on Applied Mathematics, 30 (1976), pp. 143–148.
- [43] G. MEANTI, L. CARRATINO, L. ROSASCO, AND A. RUDI, *Kernel methods through the roof: handling billions of points efficiently*, in To appear at NeurIPS 2020, preprint arXiv:2006.10350, 2020.
- [44] H. MINH, *Some Properties of Gaussian Reproducing Kernel Hilbert Spaces and Their Implications for Function Approximation and Learning Theory*, Constructive Approximation, 32 (2010), pp. 307–338.
- [45] C. MOUAT, *Fast algorithms and preconditioning techniques for fitting radial basis functions*, PhD thesis, Mathematics and Statistics, University of Canterbury, 2001.
- [46] C. MUSCO AND C. MUSCO, *Recursive Sampling for the Nyström Method*, in Advances in Neural Information Processing Systems 30, 2017, pp. 3833–3845.
- [47] M. MUTNÝ, M. DEREZIŃSKI, AND A. KRAUSE, *Convergence Analysis of Block Coordinate Algorithms with Determinantal Sampling*, in Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics (AISTATS), vol. 108, PMLR, 2020, pp. 3110–3120.
- [48] A. NOSEDAL-SANCHEZ, C. STORLIE, T. LEE, AND R. CHRISTENSEN, *Reproducing Kernel Hilbert Spaces for Penalized Regression: A Tutorial*, The American Statistician, 66 (2012), pp. 50 – 60.
- [49] A. POINAS AND R. BARDENET, *On proportional volume sampling for experimental design in general spaces*, ArXiv, abs/2011.04562 (2019).
- [50] L. PRONZATO AND A. PÁZMAN, *Design of Experiments in Nonlinear Models*, vol. 212 of Lecture Notes in Statistics, Springer-Verlag, New York, 2013.
- [51] A. RUDI, D. CALANDRIELLO, L. CARRATINO, AND L. ROSASCO, *On fast leverage score sampling and optimal learning*, in Advances in Neural Information Processing Systems, 2018, pp. 5673–5683.
- [52] A. RUDI, R. CAMORIANO, AND L. ROSASCO, *Less is more: Nyström computational regularization*, in Advances in Neural Information Processing Systems, 2015, pp. 1657–1665.
- [53] D. RUPPERT, M. WAND, AND R. CARROLL, *Semiparametric regression*, no. 12, Cambridge university press, 2003.
- [54] J. SCHREURS, M. FANUEL, AND J. SUYKENS, *Ensemble Kernel Methods, Implicit Regularization and Determinantal Point Processes*, ICML 2020 workshop on Negative Dependence and Submodularity, PMLR 119, (2020).
- [55] A. SMOLA, T.-T. FRIESS, AND B. SCHÖLKOPF, *Semiparametric support vector and linear programming machines*, in Advances in neural information processing systems, 1999, pp. 585–591.
- [56] J. SUYKENS, T. V. GESTEL, J. DE BRABANTER, B. DE MOOR, AND J. VANDEWALLE, *Least Squares Support Vector Machines*, World Scientific, Singapore, 2002.
- [57] N. TREMBLAY, S. BARTHELMÉ, AND P. AMBLARD, *Determinantal point processes for coresets*, Journal of Machine Learning Research, 20 (2019), pp. 1–70.
- [58] H. WENDLAND, *Computational aspects of radial basis function approximation*, in Topics in Multivariate Approximation and Interpolation, K. Jetter, M. Buhmann, W. Haussmann, R. Schaback, and J. Stöckler, eds., vol. 12 of Studies in Computational Mathematics, Elsevier, 2006, pp. 231 – 256.
- [59] C. WILLIAMS AND M. SEEGER, *Using the Nyström method to speed up kernel machines*, in Advances in neural information processing systems, 2001, pp. 682–688.