# Column Subset Selection and Nyström Approximation via Continuous Optimization

**Anant Mathur** [1]  **Sarat Moka** [1]  **Zdravko Botev** [1]

## Abstract

We propose a continuous optimization algorithm for the Column Subset Selection Problem (CSSP) and Nyström approximation. The CSSP and Nyström method construct low-rank approximations of matrices based on a predetermined subset of columns. It is well known that choosing the best column subset of size $k$ is a difficult combinatorial problem. In this work, we show how one can approximate the optimal solution by defining a penalized continuous loss function which is minimized via stochastic gradient descent. We show that the gradients of this loss function can be estimated efficiently using matrix-vector products with a data matrix $\mathbf{X}$ in the case of the CSSP or a kernel matrix $\mathbf{K}$ in the case of the Nyström approximation. We provide numerical results for a number of real datasets showing that this continuous optimization is competitive against existing methods.

## 1. Introduction

Recent advances in the technological ability to capture and collect data have meant that high-dimensional datasets are now ubiquitous in the fields of engineering, economics, finance, biology, and health sciences to name a few. In the case where the data collected is not labeled it is often desirable to obtain an accurate low-rank approximation for the data which is relatively low-cost to obtain and memory efficient. Such an approximation is useful to speed up downstream matrix computations that are often required in large-scale learning algorithms. The Column Subset Selection Problem (CSSP) and Nyström method are two such tools that generate low-rank approximations based on a subset of data instances or features from the dataset. The chosen subset of instances or features are commonly referred to as "landmark" points. The choice of landmark points determines how accurate the low-rank approximation is.

The challenge in the CSSP is to select the best $k$ columns of a data matrix $\mathbf{X} \in \mathbb{R}^{m \times n}$ that span its column space. That

is, for any binary vector $\boldsymbol{s} \in \{0,1\}^n$, compute

$$\operatorname*{argmin}_{\boldsymbol{s} \in \{0,1\}^n} \|\mathbf{X} - \mathbf{P}_s \mathbf{X}\|_F^2, \quad \text{subject to } \|\boldsymbol{s}\|_0 \leq k, \quad (1)$$

where $\| \cdot \|_F$ is the Frobenius matrix norm, $\|\boldsymbol{s}\|_0 = \sum_{j=1}^n I(s_j = 1)$ and $\mathbf{P}_s$ is the projection matrix onto $\text{span}\{\boldsymbol{x}_j : s_j = 1, j = 1, \dots, n\}$ ($\boldsymbol{x}_j$ being the $j$-th column of $\mathbf{X}$).

Solving this combinatorial problem exactly is known to be NP-complete (Shitov, 2021), and is practically infeasible even when $k$ is of moderate size. We propose a novel continuous optimization algorithm to approximate the exact solution to this problem. While an optimization approach via Group Lasso (Yuan & Lin, 2006) exists for the convex relaxation of this problem (Bien et al., 2010), to the best of our knowledge, no continuous optimization method has been developed to solve the highly non-convex combinatorial problem (1). To introduce our approach for the CSSP, instead of searching over binary vectors $\boldsymbol{s} \in \{0,1\}^n$, we consider the hyper-cube $[0,1]^n$ and define for each $\boldsymbol{t} \in [0,1]^n$ a matrix $\widetilde{\mathbf{P}}(\boldsymbol{t})$ which allows the following well-defined penalized continuous extension of the exact problem,

$$\operatorname*{argmin}_{\boldsymbol{t} \in [0,1]^n} \|\mathbf{X} - \widetilde{\mathbf{P}}(\boldsymbol{t})\mathbf{X}\|_F^2 + \lambda \sum_{j=1}^n t_j.$$

The parameter $\lambda > 0$ plays an analogous role to that of the regularization parameter in regularized linear regression methods (Tibshirani, 1996) and controls the sparsity of the solution, that is, the size of $k$. Two aspects of this continuous extension make it useful for approximating the exact solution. Firstly, the continuous loss agrees with the discrete loss at every corner point $\boldsymbol{s} \in \{0,1\}^n$ of the hypercube $[0,1]^n$, and secondly, for large datasets the gradient can be estimated via an unbiased stochastic estimate. To obtain an approximate solution to the exact problem, *stochastic gradient descent* (SGD) is implemented on the penalized loss. After starting at an interior point of the hyper-cube, under SGD, the vector $\boldsymbol{t}$ moves towards a corner point, and some of the $t_j$'s exhibit shrinkage to zero. It is these values that indicate which columns in $\mathbf{X}$ should not be selected as landmark points.

The Nyström approximation (Williams & Seeger, 2000; Drineas et al., 2005) is a popular variant of the CSSP for

positive semi-definite kernel matrices. The Nyström method also constructs a low-rank approximation $\widehat{\mathbf{K}} \in \mathbb{R}^{n \times n}$ to the true kernel matrix $\mathbf{K} \in \mathbb{R}^{n \times n}$ using a subset of columns. Once the $k$ columns are selected, $\widehat{\mathbf{K}}$ (in factored form) takes $O(k^3)$ additional time to compute, requires $O(nk)$ space to store, and can be manipulated quickly in downstream applications, e.g., inverting $\widehat{\mathbf{K}}$ takes $O(nk^2)$ time. In addition to the continuous extension for the CSSP, in this paper, we provide a continuous optimization algorithm that can approximate the best $k$ columns to be used to construct $\widehat{\mathbf{K}}$ (Section 2.2).

The continuous algorithm for the CSSP formulated in this paper utilizes SGD where at each iteration one can estimate the gradient with a cost of $O(mn)$. We show that the gradients of the penalized continuous loss can be estimated via linear solves with random vectors that are approximated with the conjugate gradient algorithm (CG) (Golub & Van Loan, 1996), which itself is an iterative algorithm that only requires matrix-vector multiplications (MVMs) with the $m \times n$ matrix $\mathbf{X}$. Similarly, for the Nyström method we show that at each step of the gradient descent, the gradient can be estimated in $O(n^2)$ time requiring only matrix-vector multiplications with the kernel matrix $\mathbf{K}$. This is especially useful in cases where we only have access to a black-box MVM function. The fact that both these algorithms require only matrix-vector multiplications to estimate the gradients lends itself to utilizing GPU hardware acceleration. Moreover, the computations in the proposed algorithm can exploit the sparsity that is achieved by working only with the columns of $\mathbf{X}$ that are selected by the algorithm at any given iteration.

### 1.1. Related Work

There exists extensive literature on random sampling methods for the approximation of the exact CSSP and Nyström problem. Sampling techniques such as adaptive sampling (Deshpande & Vempala, 2006), ridge leverage scores (Gittens & Mahoney, 2013; Musco & Musco, 2017; Alaoui & Mahoney, 2015) attempt to sample "important" and "diverse" columns. In particular, recent attention has been paid to Determinantal Point Processes (DPPs) (Hough et al., 2006; Derezinski & Mahoney, 2021). DPPs provide strong theoretical guarantees (Derezinski et al., 2020) for the CSSP and Nyström approximation and are amenable to efficient numerical implementation (Li et al., 2016; Derezinski et al., 2019; Calandriello et al., 2020; Dereziński, 2019). Outside of sampling methods, iterative methods such as Greedy selection (Farahat et al., 2011; 2013) have been shown to perform well in practice and exhibit provable guarantees (Altschuler et al., 2016).

Column selection has been extensively studied in the supervised context of linear regression (more commonly referred to as feature or variable selection). Penalized regression methods such as the Lasso (Tibshirani, 1996) have been widely applied to select columns of a predictor matrix that best explain a response vector. The canonical $k$-best subset or $l_0$-penalized regression problem is another penalized regression method, where the goal is to find the best subset of $k$ predictors that best fit a response $\boldsymbol{y}$ (Beale et al., 1967; Hocking & Leslie, 1967). The recently proposed *Continuous Optimization Method Towards Best Subset Selection* (COMBSS) algorithm (Moka et al., 2022) attempts to solve the $l_0$-penalized regression problem by minimizing a continuous loss that approximates the exact solution. The algorithm we propose for the CSSP in this paper can be viewed as an adaptation of COMBSS to the unsupervised setting. In this setting, the goal is to find the best subset of size $k$ for a multiple multivariate regression model where both the response and predictor matrix are $\mathbf{X}$. Interestingly, this framework can be extended to include a continuous selection loss for the Nyström approximation.

The rest of the paper is structured as follows. In Section 2 we describe the continuous extension for the CSSP and the Nyström method. In Section 3 we provide steps for the efficient implementation of our proposed continuous algorithm on large matrices and in Section 4 we provide numerical results on a variety of real datasets.

## 2. Continuous Loss for Landmark Selection

In this section, we formally define the CSSP and the best size $k$-Nyström approximation. Then, we provide the mathematical setup for the continuous extension of the exact problem.

### 2.1. Column Subset Selection

Let $\mathbf{X} \in \mathbb{R}^{m \times n}$ and for any binary vector $\boldsymbol{s} = (s_1, \ldots, s_n)^\top \in \{0, 1\}^n$, let $\mathbf{X}_{[\boldsymbol{s}]}$ denote the matrix of size $m \times \|s\|_0$ keeping only columns $j$ of $\mathbf{X}$ where $s_j = 1$, for $j = 1, \ldots, n$. Then for every integer $k \leq n$ the CSSP finds

$$\operatorname*{argmin}_{\boldsymbol{s} \in \{0,1\}^n} \|\mathbf{X} - \mathbf{P}_s \mathbf{X}\|_F^2, \quad \text{subject to } \|\boldsymbol{s}\|_0 \leq k, \quad (2)$$

where $\mathbf{P}_s := \mathbf{X}_{[\boldsymbol{s}]} \mathbf{X}_{[\boldsymbol{s}]}^\dagger$ ($\dagger$ denotes Moore–Penrose inverse) is the projection matrix onto $\operatorname{span}\{\mathbf{x}_j : s_j = 1\}$ and $\boldsymbol{x}_j$ is the $j$-th column of $\mathbf{X}$. By expanding the Frobenius norm it is easy to see that the discrete problem (2) can be reformulated as,

$$\operatorname*{argmin}_{\boldsymbol{s} \in \{0,1\}^n} -\operatorname{tr}\left[\mathbf{X}^\top \mathbf{P}_{\boldsymbol{s}} \mathbf{X}\right], \quad \text{subject to } \|\boldsymbol{s}\|_0 \leq k.$$

We now define a new matrix function on $\boldsymbol{t} \in [0, 1]^n$ which acts as a continuous generalization of $\mathbf{P}_{\boldsymbol{s}}$.

**Definition 2.1.** For $\boldsymbol{t} = (t_1, \ldots, t_n)^\top \in [0,1]^n$, define $\mathbf{T} := \mathrm{Diag}(\boldsymbol{t})$ as the diagonal matrix with diagonal elements $t_1, \ldots, t_n$ and

$$\widetilde{\mathbf{P}}(\boldsymbol{t}) := \mathbf{XT} \left[ \mathbf{TX}^\top \mathbf{XT} + \delta(\mathbf{I} - \mathbf{T}^2) \right]^\dagger \mathbf{TX}^\top,$$

where $\delta > 0$ is a fixed constant.

Although not explicitly stated in (Moka et al., 2022), $\widetilde{\mathbf{P}}(\boldsymbol{t})$ is used as the continuous generalization for the hat matrix $\mathbf{P}_{\boldsymbol{s}}$ to solve the $l_0$-penalized regression problem.

The main difference between this definition and traditional sampling methods is that instead of multiplying $\mathbf{X}$ by a sampling matrix to obtain $\mathbf{X}_{[\boldsymbol{s}]}$ we compute the matrix $\mathbf{XT}$ which weights column $j$ of $\mathbf{X}$ by the parameter $t_j \in [0,1]$. Intuitively, the matrix $\mathbf{TX}^\top \mathbf{XT} + \delta(\mathbf{I} - \mathbf{T}^2)$ can be viewed as a convex combination of the matrices $\mathbf{X}^\top \mathbf{X}$ and $\delta \mathbf{I}$.

From an evaluation standpoint, the pseudo-inverse need not be evaluated for any interior point in this newly defined function. We remark that for any $\boldsymbol{t} \in [0,1)^n$ the matrix inverse in Definition 2.1 exists and therefore,

$$\widetilde{\mathbf{P}}(\boldsymbol{t}) = \mathbf{XT} \left[ \mathbf{TX}^\top \mathbf{XT} + \delta(\mathbf{I} - \mathbf{T}^2) \right]^{-1} \mathbf{TX}^\top.$$

We now state two results for the function $\widetilde{\mathbf{P}}(\boldsymbol{t})$ and its relationship with the projection matrix $\mathbf{P}_s$. The following Lemmas (2.2 and 2.3) are extensions of the results stated in (Moka et al., 2022).

**Lemma 2.2.** *For any binary vector $\boldsymbol{s} \in \{0,1\}^n$, $\widetilde{\mathbf{P}}(\boldsymbol{s})$ exists and*

$$\widetilde{\mathbf{P}}(\boldsymbol{s}) = \mathbf{P}_{\boldsymbol{s}} = \mathbf{X}_{[\boldsymbol{s}]} \mathbf{X}_{[\boldsymbol{s}]}^\dagger.$$

**Lemma 2.3.** *$\widetilde{\mathbf{P}}(\boldsymbol{t})$ is continuous element-wise over $[0,1]^n$. Moreover, for any sequence $\boldsymbol{t}^{(1)}, \boldsymbol{t}^{(2)} \cdots \in [0,1)^n$ converging to $\boldsymbol{t} \in [0,1]^n$, the limit $\lim_{l \to \infty} \widetilde{\mathbf{P}}(\boldsymbol{t}^{(l)})$ exists and*

$$\lim_{l \to \infty} \widetilde{\mathbf{P}}(\boldsymbol{t}^{(l)}) = \widetilde{\mathbf{P}}(\boldsymbol{t}).$$

We note that the proof of Lemma 2.3 follows identically to the proof of *Theorem 3* in (Moka et al., 2022) where it is stated that the function $\|\boldsymbol{y} - \widetilde{\mathbf{P}}(\boldsymbol{t})\boldsymbol{y}\|_2^2$ is continuous over $[0,1]^n$ for any fixed vector $\boldsymbol{y} \in \mathbb{R}^n$.

Given $\widetilde{\mathbf{P}}(\boldsymbol{t})$ is continuous on $[0,1]^n$ and agrees with $\mathbf{P}_{\boldsymbol{s}}$ at every corner point we can define the continuous generalization of the exact problem (2),

$$\underset{\boldsymbol{t} \in [0,1]^n}{\mathrm{argmin}} - \mathrm{tr} \left[ \mathbf{X}^\top \widetilde{\mathbf{P}}(\boldsymbol{t}) \mathbf{X} \right], \quad \text{subject to } \sum_{j=1}^{n} t_j \leq k.$$

Instead of solving this constrained problem, for a tunable parameter $\lambda$, we consider minimizing the Lagrangian function,

$$\underset{\boldsymbol{t} \in [0,1]^n}{\mathrm{argmin}} f_\lambda(\boldsymbol{t}), \quad f_\lambda(\boldsymbol{t}) := - \mathrm{tr} \left[ \mathbf{X}^\top \widetilde{\mathbf{P}}(\boldsymbol{t}) \mathbf{X} \right] + \lambda \sum_{j=1}^{n} t_j.$$

In Section 3 we reformulate this box-constrained problem into an equivalent unconstrained problem via a nonlinear mapping $\boldsymbol{t} = \boldsymbol{t}(\boldsymbol{w})$ for $\boldsymbol{w} \in \mathbb{R}^n$ that forces $\boldsymbol{t}$ to be in the hypercube $[0,1]^n$. We solve this optimization via continuous gradient descent. To this end, we need to evaluate the gradient $\nabla f_\lambda(\boldsymbol{t})$ for any interior point.

**Lemma 2.4.** *Let $\mathbf{K} = \mathbf{X}^\top \mathbf{X}$, $\mathbf{Z} = \mathbf{K} - \delta \mathbf{I}$ and $\mathbf{L}_{\boldsymbol{t}} = \mathbf{TZT} + \delta \mathbf{I}$. Then, for $\boldsymbol{t} \in (0,1)^n$,*

$$\nabla f_\lambda(\boldsymbol{t}) = 2 \, \mathrm{Diag} \left[ \mathbf{L}_{\boldsymbol{t}}^{-1} \mathbf{TK}^2 \left( \mathbf{TL}_{\boldsymbol{t}}^{-1} \mathbf{TZ} - \mathbf{I} \right) \right] + \lambda \mathbf{1}.$$

Evaluating $\nabla f_\lambda(\boldsymbol{t})$ has a computational complexity of $O(n^3)$ due to the required inversion of $\mathbf{L}_t$. In Section 3 we detail an unbiased estimate for $\nabla f_\lambda(\boldsymbol{t})$ which utilizes the CG algorithm, where the most expensive operations involved are matrix-vector multiplications with $\mathbf{X}$ and $\mathbf{X}^\top$, which reduces the computational complexity to $O(mn)$.

## 2.2. Nyström Method

We now turn our attention to defining a continuous objective for the landmark points in the Nyström approximation. We consider optimizing the landmark points first with respect to the trace matrix norm and then to the Frobenius matrix norm.

In many applications, we are interested in obtaining a low-rank approximation to a *kernel* matrix $\mathbf{K} \in \mathbb{R}^{n \times n}$. Consider an input space $\mathcal{X}$ and a positive semi-definite kernel function $h : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$. Given a set of $n$ input points $\boldsymbol{x}_1', ..., \boldsymbol{x}_n' \in \mathcal{X}$, the kernel matrix $\mathbf{K} \in \mathbb{R}^{n \times n}$ is defined by $\mathbf{K}_{i,j} = h(\boldsymbol{x}_i', \boldsymbol{x}_j')$ and is positive semi-definite.

For any binary vector $\boldsymbol{s} \in \{0,1\}^n$ let $\mathbf{K}_{[\boldsymbol{s}]}$ be the $n \times \|\boldsymbol{s}\|_0$ matrix with columns indexed by $\{j : s_j = 1\}$ and $\mathbf{K}_{[\boldsymbol{s},\boldsymbol{s}]}$ be the $\|\boldsymbol{s}\|_0 \times \|\boldsymbol{s}\|_0$ principal sub-matrix indexed by $\{j : s_j = 1\}$. The Nyström low-rank approximation for $\mathbf{K}$ is given by,

$$\widehat{\mathbf{K}}_s := \mathbf{K}_{[\boldsymbol{s}]} \mathbf{K}_{[\boldsymbol{s},\boldsymbol{s}]}^\dagger \mathbf{K}_{[\boldsymbol{s}]}^\top.$$

The following observation appearing in (Derezinski et al., 2020) connects the CSSP and the Nyström approximation with respect to the trace matrix norm.

Suppose we have the decomposition of the kernel matrix $\mathbf{K} = \mathbf{X}^\top \mathbf{X}$ where $\mathbf{X} \in \mathbb{R}^{m \times n}$. Then, the Nyström approximation is given by $\widehat{\mathbf{K}}_{\boldsymbol{s}} = (\mathbf{P}_{\boldsymbol{s}} \mathbf{X})^\top \mathbf{P}_{\boldsymbol{s}} \mathbf{X}$ and

$$\|\mathbf{K} - \widehat{\mathbf{K}}_{\boldsymbol{s}}\|_* = \|\mathbf{X} - \mathbf{P}_{\boldsymbol{s}} \mathbf{X}\|_F^2.$$

where $\|\mathbf{A}\|_* = \sum_{i=1}^{\min\{m,n\}} \sigma_i(\mathbf{A})$ for $\mathbf{A} \in \mathbb{R}^{m \times n}$ is the trace matrix norm. This connection is used in (Derezinski et al., 2020) to provide shared approximation bounds for both the CSSP and Nyström approximation. Given that the kernel matrix is always positive semi-definite, the decomposition $\mathbf{K} = \mathbf{X}^\top \mathbf{X}$ always exists and one can solve the

CSSP for $\mathbf{X}$ to obtain the best $k$-landmark Nyström approximation with respect to the trace norm. We note that such a decomposition is not unique, e.g., it can be the Cholesky decomposition or the symmetric square-root decomposition.

The matrix $\mathbf{X}$ does not need explicit evaluation in order to perform CSSP as one can attain $\nabla f_\lambda(\boldsymbol{t})$ with the matrix $\mathbf{K}$ instead (see, Lemma 2.4). Therefore, finding the decomposition $\mathbf{K} = \mathbf{X}^\top \mathbf{X}$ is not required, and one can approximately solve the CSSP by minimizing $\nabla f_\lambda(\boldsymbol{t})$ with the kernel matrix $\mathbf{K}$.

Suppose instead we want to use the Frobenius matrix norm to find the best choice of columns of the matrix $\mathbf{K}$ to construct the Nyström approximation. This problem is formulated as

$$\operatorname*{argmin}_{\boldsymbol{s} \in \{0,1\}^n} \|\mathbf{K} - \widehat{\mathbf{K}}_{\boldsymbol{s}}\|_F^2, \quad \text{subject to } \|\boldsymbol{s}\|_0 \le k. \quad (3)$$

Similar to $\widetilde{\mathbf{P}}(\boldsymbol{t})$ we can weight each column $j$ of $\mathbf{K}$ by $t_j \in [0,1]$ instead of sampling the columns $\mathbf{K}_{[\boldsymbol{s}]}$ for the Nyström approximation. We define continuous generalization for the Nyström approximation,

**Definition 2.5.** For $\boldsymbol{t} = (t_1, \ldots, t_n)^\top \in [0,1]^n$ let $\mathbf{T} := \operatorname{Diag}(\boldsymbol{t})$ and

$$\widetilde{\mathbf{K}}(\boldsymbol{t}) := \mathbf{K}\mathbf{T}\left[\mathbf{T}\mathbf{K}\mathbf{T} + \delta(\mathbf{I} - \mathbf{T}^2)\right]^\dagger \mathbf{T}\mathbf{K},$$

where $\delta > 0$ is a fixed constant. Similar to $\widetilde{\mathbf{P}}(\boldsymbol{t})$, for any $\boldsymbol{t} \in [0,1)^n$ the matrix $\mathbf{T}\mathbf{K}\mathbf{T} + \delta(\mathbf{I} - \mathbf{T}^2)$ is invertible.

In the following two results, we state that $\widetilde{\mathbf{K}}(\boldsymbol{t})$ is a continuous function on $[0,1]^n$ and agrees with the exact Nyström approximation at every corner point.

**Lemma 2.6.** For any corner point $\boldsymbol{s} \in \{0,1\}^n$, $\widetilde{\mathbf{K}}(\boldsymbol{s})$ exists and

$$\widetilde{\mathbf{K}}(\boldsymbol{s}) = \widehat{\mathbf{K}}_s = \mathbf{K}_{[\boldsymbol{s}]}\mathbf{K}_{[\boldsymbol{s},\boldsymbol{s}]}^\dagger \mathbf{K}_{[\boldsymbol{s}]}^\top.$$

**Lemma 2.7.** $\widetilde{\mathbf{K}}(\boldsymbol{t})$ is continuous element-wise over $[0,1]^n$. Moreover, for any sequence $\boldsymbol{t}^{(1)}, \boldsymbol{t}^{(2)} \cdots \in [0,1)^n$ converging to $\boldsymbol{t} \in [0,1]^n$, the limit $\lim_{l \to \infty} \widetilde{\mathbf{K}}(\boldsymbol{t}^{(l)})$ exists and

$$\lim_{l \to \infty} \widetilde{\mathbf{K}}(\boldsymbol{t}^{(l)}) = \widetilde{\mathbf{K}}(\boldsymbol{t}).$$

We therefore have the continuous generalization of the exact problem (3),

$$\operatorname*{argmin}_{\boldsymbol{t} \in [0,1]^n} \|\mathbf{K} - \widetilde{\mathbf{K}}(\boldsymbol{t})\|_F^2, \quad \text{subject to } \sum_{j=1}^n t_j \le k.$$

Instead of solving this constrained problem, for a tunable parameter $\lambda$, we consider minimizing the Lagrangian function,

$$\operatorname*{argmin}_{\boldsymbol{t} \in [0,1]^n} g_\lambda(\boldsymbol{t}), \quad g_\lambda(\boldsymbol{t}) := \|\mathbf{K} - \widetilde{\mathbf{K}}(\boldsymbol{t})\|_F^2 + \lambda \sum_{j=1}^n t_j.$$

As with the continuous extension for CSSP we use a gradient descent method to solve the above problem. The following result provides an expression for $\nabla g_\lambda(\boldsymbol{t})$ for $\boldsymbol{t} \in (0,1)^n$.

**Lemma 2.8.** Let $\mathbf{Z} = \mathbf{K} - \delta\mathbf{I}$, $\mathbf{L}_{\boldsymbol{t}} = \mathbf{T}\mathbf{Z}\mathbf{T} + \delta\mathbf{I}$ and $\mathbf{D} = \widetilde{\mathbf{K}}(\boldsymbol{t}) - \mathbf{K}$. Then, for $\boldsymbol{t} \in (0,1)^n$,

$$\nabla g_\lambda(\boldsymbol{t}) = 4\operatorname{Diag}\left[\mathbf{L}_{\boldsymbol{t}}^{-1}\mathbf{T}\mathbf{K}\mathbf{D}\mathbf{K}\left(\mathbf{I} - \mathbf{T}\mathbf{L}_{\boldsymbol{t}}^{-1}\mathbf{T}\mathbf{Z}\right)\right] + \lambda\mathbf{1}.$$

Evaluating $\nabla g_\lambda(\boldsymbol{t})$ has a computational complexity of $O(n^3)$ due to the required inversion of $\mathbf{L}$ and evaluation of $\mathbf{K}(\boldsymbol{t})$. As with $\nabla f_\lambda(\boldsymbol{t})$ we detail an unbiased estimate for $\nabla g_\lambda(\boldsymbol{t})$ in Section 3 which utilizes matrix-vector multiplications with $\mathbf{K}$ and that helps in reducing the computational cost.

# 3. Implementation

In this section, we detail how to efficiently solve the continuous problems posed in Section 2. In particular, we detail a non-linear transformation that was also used in (Moka et al., 2022) to make both the CSSP and Nyström approximation optimization problems unconstrained. We then show how one can estimate the gradients using MVMs with $\mathbf{X}$ and $\mathbf{K}$.

### 3.1. Handling Box Constraints (Moka et al., 2022)

The continuous extension of the CSSP and Nyström approximation requires minimizing the functions $f_\lambda(\boldsymbol{t})$ and $g_\lambda(\boldsymbol{t})$ over $\boldsymbol{t} \in [0,1]^n$. We now consider a non-linear transformation to make both optimization problems unconstrained. Consider the mapping $\boldsymbol{t} = \boldsymbol{t}(\boldsymbol{w})$ given by,

$$t_j(w_j) = 1 - \exp(-w_j^2), \quad j = 1, \ldots, n,$$

then if we consider the optimization of continuous CSSP,

$$\boldsymbol{w}^* = \operatorname*{argmin}_{\boldsymbol{w} \in \mathbb{R}^p} f_\lambda(\boldsymbol{t}(\boldsymbol{w})),$$

we attain the solution to (2.1) by evaluating $\boldsymbol{t}(\boldsymbol{w}^*)$. This is true because for any $a, b \in \mathbb{R}$,

$$1 - \exp(-a^2) < 1 - \exp(-b^2) \quad \text{if and only if } a^2 < b^2.$$

In vector form the transformation is $\boldsymbol{t}(\boldsymbol{w}) = \mathbf{1} - \exp(-\boldsymbol{w} \odot \boldsymbol{w})$ (here $\odot$ denotes element-wise multiplication) and using the chain rule we obtain for $\boldsymbol{w} \in \mathbb{R}^p$,

$$\frac{\partial f_\lambda(\boldsymbol{t}(\boldsymbol{w}))}{\partial \boldsymbol{w}} = \frac{\partial f_\lambda(\boldsymbol{t}(\boldsymbol{w}))}{\partial \boldsymbol{t}} \odot (2\boldsymbol{w} \odot \exp(-\boldsymbol{w} \odot \boldsymbol{w})).$$

We can now implement a gradient descent algorithm to approximately obtain $\boldsymbol{t}(\boldsymbol{w}^*)$. Using this approximation we can select an appropriate binary vector as a solution to the exact problem (2). The same transformation can be applied to solve $g_\lambda(\boldsymbol{t}(\boldsymbol{w}))$ over $\boldsymbol{w} \in \mathbb{R}^n$.

## 3.2. Stochastic Estimate for the Gradient

As discussed in Section 2, $\nabla f_\lambda(t)$ and $\nabla g_\lambda(t)$ are problematic to compute for large $n$ due to the $O(n^3)$ complexity of inverting a matrix. Here we show that we can implement a stochastic gradient descent (SGD) which has strong theoretical guarantees (Robbins & Monro, 1951) by using an unbiased estimate for $\nabla f_\lambda(t)$ and $\nabla g_\lambda(t)$.

The method we employ is a factorized estimator $\hat{\ell}$ for the diagonal of a square matrix. Suppose we wish to estimate the diagonal of the matrix $\mathbf{A} = \mathbf{B}\mathbf{C}^\top$ where $\mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathbb{R}^{n \times n}$. Let $z \in \mathbb{R}^n$ be a random vector sampled from the Rademacher distribution, whose entries are either $-1$ or $1$, each with probability $1/2$. Then an unbiased estimate for $\mathrm{Diag}(\mathbf{A})$ is $\hat{\ell} = \mathbf{B}z \odot \mathbf{C}z$, see (Martens et al., 2012). Further analysis of its properties including its variance can be found in (Mathur et al., 2021). We note that when $\mathbf{B} = \mathbf{A}$ and $\mathbf{C} = \mathbf{I}$, this estimator reduces to the well-known (Bekas et al., 2007) estimator for the diagonal.

The two following results provide an unbiased estimate for $\nabla f_\lambda(t)$ and $\nabla g_\lambda(t)$ using the factorized estimator for the diagonal of a matrix.

**Lemma 3.1.** *Recall that in the continuous CSSP optimization for* $\mathbf{X}$*, we have the definitions* $\mathbf{T} = \mathrm{Diag}(t)$ *for* $t \in [0,1]^n$*,* $\mathbf{K} = \mathbf{X}^\top \mathbf{X}$*,* $\mathbf{Z} = \mathbf{K} - \delta \mathbf{I}_n$ *and* $\mathbf{L}_t = \mathbf{T}\mathbf{Z}\mathbf{T} + \delta \mathbf{I}_n$*.*

*Suppose* $z \in \mathbb{R}^n$ *follows a Rademacher distribution and let:*

*(1)* $a = \mathbf{K}z$*, (2)* $b = \mathbf{L}_t^{-1}(t \odot a)$ *and*

$$\phi = b \odot \mathbf{Z}(t \odot b) - a \odot b.$$

*Then for* $t \in (0,1)^n$*,*

$$\nabla f_\lambda(t) = 2\mathbb{E}[\phi] + \lambda \mathbf{1}.$$

**Lemma 3.2.** *Recall that in the continuous Nyström optimization for a kernel matrix* $\mathbf{K}$*, we have the definitions* $\mathbf{T} = \mathrm{Diag}(t)$ *for* $t \in [0,1]^n$*,* $\mathbf{Z} = \mathbf{K} - \delta \mathbf{I}$ *and* $\mathbf{L}_t = \mathbf{T}\mathbf{Z}\mathbf{T} + \delta \mathbf{I}$*.*

*Suppose* $z \in \mathbb{R}^n$ *follows a Rademacher distribution and let:*

*(1)* $a = \mathbf{K}z$*, (2)* $b = \mathbf{L}_t^{-1}(t \odot a)$*, (3)* $c = \mathbf{K}(t \odot b) - a$*, (4)* $d = \mathbf{K}c$*, (5)* $e = \mathbf{L}_t^{-1}(t \odot d)$ *and*

$$\psi = b \odot d + a \odot e - e \odot \mathbf{Z}(t \odot b) - b \odot \mathbf{Z}(t \odot e).$$

*Then for* $t \in (0,1)^n$*,*

$$\nabla g_\lambda(t) = 2\mathbb{E}[\psi] + \lambda \mathbf{1}.$$

Using these results, we can obtain for a Monte-Carlo size $M$, the approximations $\nabla f_\lambda(t) \approx 2\left(\frac{1}{M}\sum_{i=1}^M \phi^{(i)}\right) + \lambda \mathbf{1}$ and $\nabla g_\lambda(t) \approx 2\left(\frac{1}{M}\sum_{i=1}^M \psi^{(i)}\right) + \lambda \mathbf{1}$, where $\phi^{(i)}$ and $\psi^{(i)}$ are evaluated using a sample $z^{(i)}$ drawn from the Rademacher distribution.

---

**Algorithm 1** Continous Landmark Selection

---

1: **input:** Data matrix: $\mathbf{X} \in \mathbb{R}^{m \times n}$ (CSSP) or Kernel matrix: $\mathbf{K} \in \mathbb{R}^{n \times n}$ (Nyström method), Tuning parameters: $\delta$ and $\lambda$, Monte Carlo size: $M$, Termination Condition: TermCond, Threshold value: $\tau \in [0,1]$.
2: Set $t^{(0)} = (1/2, \dots 1/2)^\top$
3: $w^{(0)} \leftarrow \sqrt{-\ln(1 - t^{(0)})}$
4: $w^* \leftarrow$ SGD $(w^{(0)}, M, \mathbf{X}$ or $\mathbf{K}$, TermCond$)$
5: $t^* \leftarrow 1 - \exp(-w^* \odot w^*)$
6: **for** $i = 1$ **to** $n$ **do**
7:   $s_j \leftarrow I(t_j^* > \tau)$
8: **end for**
9: **return:** $s^* = (s_1, \dots, s_n)^\top$

---

These results show that to evaluate stochastic gradients one needs to solve linear systems efficiently with the matrix $\mathbf{L}_t$. These systems can be iteratively solved using the conjugate gradient (CG) algorithm (Golub & Van Loan, 1996) which uses a sequence of MVMs with $\mathbf{L}_t$. Multiplying a vector with $\mathbf{L}_t$ can be reduced to a single MVM with the matrix $\mathbf{K}$ and a sequence of element-wise vector multiplications and additions.

### 3.3. Obtaining a Solution

While we have re-framed both the CSSP and the Nyström problem as an optimization over $t \in [0,1]^n$, the priority remains to obtain an approximate solution $s \in \{0,1\}^n$ to (2) and (3). To obtain such a binary vector, we first initialize SGD from a starting point $t^{(0)}$ and return the final value $t^*$ after a termination condition for SGD has been satisfied. Under SGD the iterative sequence $\{t^{(i)}\}_{i\geq 0}$ moves towards a corner point of the hypercube. To obtain the closest corner point $s \in \{0,1\}^n$, we map the insignificant $t_j^*$'s to 0 and all the other $t_j^*$'s to 1 for some tolerance parameter $\tau \in (0,1)$. This implementation is shown Algorithm 1. In Figure 1 we provide example solution paths $\{t^{(i)}\}_{i\geq 0}$ under both batch gradient descent and SGD.

When choosing the value for $t^{(0)}$ it is important to consider the following true statements: $t_j = 0$ if and only if $w_j = 0$ and

$$\lim_{w_j \to 0} \frac{\partial f_\lambda(t(w))}{\partial w_j} = \lim_{w_j \to 0} \frac{\partial g_\lambda(t(w))}{\partial w_j} = 0.$$

These facts imply that if $t_j$ is set to zero during the course of the optimization it will remain unchanged thereafter. Therefore, it is important to choose $t^{(0)}$ that is away from any corner point. It is for this reason, we set $t^{(0)} = (1/2, \dots, 1/2)^\top$ in all our experiments.

(a) Gradient Descent

(b) Stochastic Gradient Descent

*Figure 1.* Convergence of $\boldsymbol{t}$ for continuous Column Subset Selection using the MNIST dataset. Blue trajectories correspond to selected columns. Only a subset of 300 randomly chosen column trajectories (out of 784) are displayed. For both (a) and (b), $\lambda = 10$ and $\delta = 10$. In (b) the Monte-Carlo size is $M = 5$.

### 3.4. Dimensionality Reduction

In Section 3.3 we stated that if $t_j$ is set to zero during the course of the SGD then it will remain unchanged thereafter. This opens the possibility to reduce the computational cost of estimating $\nabla f_\lambda(\boldsymbol{t}(\boldsymbol{w}))$ and $\nabla g_\lambda(\boldsymbol{t}(\boldsymbol{w}))$ by only focusing on terms where $t_j \neq 0$.

Let $\mathcal{N} = \{1, \ldots, n\}$ and for any $\boldsymbol{t} \in [0, 1)^n$ let $\mathcal{I}_{\boldsymbol{t}} = \{j : t_j = 0\}$. For a vector $\boldsymbol{a} \in \mathbb{R}^n$, denote the vector $(\boldsymbol{a})_+$ of dimension $n - |\mathcal{I}_{\boldsymbol{t}}|$ (respectively $|\mathcal{I}_{\boldsymbol{t}}|$) constructed from $\boldsymbol{a}$ by removing the elements with indices that are in $\mathcal{I}_{\boldsymbol{t}}$ (respectively, in $\mathcal{N} \setminus \mathcal{I}_{\boldsymbol{t}}$). Likewise, for a matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, denote the principal sub-matrix $(\mathbf{A})_+$ (respectively, $(\mathbf{A})_0$) that is constructed by removing the rows and columns with indices that are in $|\mathcal{I}_{\boldsymbol{t}}|$ (respectively, in $\mathcal{N} \setminus \mathcal{I}_{\boldsymbol{t}}$). Then, we have the following result.

**Lemma 3.3.** *For any expression $\boldsymbol{q} = \mathbf{L}_{\boldsymbol{t}}^{-1}(\boldsymbol{t} \odot \boldsymbol{r})$ where $\boldsymbol{r} \in \mathbb{R}^n$ and $\boldsymbol{t} \in [0, 1)^n$,*

$$(\boldsymbol{q})_0 = \mathbf{0} \quad and \quad (\boldsymbol{q})_+ = ((\mathbf{L}_{\boldsymbol{t}})_+)^{-1}((\boldsymbol{t})_+ \odot (\boldsymbol{r})_+),$$

*where,*

$$(\mathbf{L}_{\boldsymbol{t}})_+ = (\mathbf{T})_+(\mathbf{K})_+(\mathbf{T})_+ + \delta(\mathbf{I} - (\mathbf{T})_+^2).$$

To incorporate this result in our algorithm, we set a small constant $\epsilon$ and during the course of SGD if $\boldsymbol{t}(\boldsymbol{w})_j < \epsilon$, we set its value to zero. Thereafter, when solving (2) and (5) in either Lemma 3.1 (CSSP) or Lemma 3.2 (Nyström) the dimension of the linear system is $n - |\mathcal{I}_{\boldsymbol{t}}| < n$.

### 3.5. Complexity Analysis

The main computational cost of our algorithm is the complexity attributed to estimating the gradients at each iteration of SGD. For simplicity of analysis, we assume the dimensionality reduction described in Section 3.4 is not carried

out. The cost to solve (2) and (5) in either Lemma 3.1 or 3.2 via CG is $O(T_{mult} M\ell)$ flops where $\ell$ is the number of CG iterations and $T_{mult}$ is the cost of computing a matrix-vector product with either $\mathbf{X}^\top \mathbf{X}$ (CSSP) or kernel matrix $\mathbf{K}$ (Nyström). Generally, only $\ell \ll n$ iterations of CG are required to obtain an accurate solution to the linear system.

The cost $T_{mult}$ is $O(mn)$ and $O(n^2)$ via direct computation for $\mathbf{X}^\top \mathbf{X}$ and $\mathbf{K}$ respectively. For kernel matrices with specific structure, this cost can be reduced. For example, for Toeplitz matrices or for matrices constructed from a kernel function that is analytic and isotropic, the cost can be reduced to quasi-linear complexity (Dietrich & Newsam, 1997; Gardner et al., 2018; Ryan et al., 2022). Utilizing GPU hardware for accelerating matrix computations has gained significant recent attention and numerous software regimes (Charlier et al., 2021; Hu et al., 2022) have been proposed to accelerate kernel MVMs. These methods can be implemented out-of-the-box and allow MVMs to be feasible on very large datasets ($n \sim 10^8$). Another advantage of these algorithms is that, as long as the kernel function $h(\boldsymbol{x}'_i, \boldsymbol{x}'_j)$ is given, MVMs can be computed directly without ever storing the kernel matrix $\mathbf{K}$. This is an advantage of our method when compared to other methods such as the greedy selection method for the Nyström approximation in (Farahat et al., 2011), which has a cost of $O(n^2 k)$ and requires the full explicit matrix to be stored in memory.

### 3.6. Role of parameters $\delta$ and $\lambda$

The tuning parameter $\lambda$ controls the size of the penalty $\|\boldsymbol{t}\|_1$ which is added to the Frobenius matrix loss. It is intuitive then that for a larger value of $\lambda$ a stronger shrinkage is applied to $\boldsymbol{t}$ during the course of the continuous optimization. In terms of curvature, as $\lambda$ increases so does the directional slope of $f_\lambda(\boldsymbol{t}(\boldsymbol{w}))$ and $g_\lambda(\boldsymbol{t}(\boldsymbol{w}))$ in the region around $w_j = 0$. For this reason, it is likelier that more $w_j$'s will be pushed

*Figure 2.* The mean Nyström empirical approximation factor over 50 trials for the UCI Residential Building and MNIST dataset where $\mathbf{K}$ is constructed using the Gaussian Radial Basis Function (RBF) kernel: $\mathbf{K}_{i,j} = h(\boldsymbol{x}'_i, \boldsymbol{x}'_j) = \exp\left(-\|\boldsymbol{x}'_i - \boldsymbol{x}'_j\|^2\right)/\sigma^2$. Approximation factor is plotted on a logarithmic scale.



*Figure 3.* The mean CSSP empirical approximation factor over 50 trials for the MNIST dataset and three UCI datasets for different methods. Approximation factor is plotted on a logarithmic scale.

towards zero when the value for $\lambda$ is large. This behavior is similar to that of the parameter $\lambda$ in the COMBSS method (Moka et al., 2022) where a more formal analysis can be found. We note that the relationship between $\lambda$ and $k$ is data dependent and it is suggested that the user apply an efficient grid search regime to obtain an appropriate $\lambda$ for their use.

With respect to the parameter $\delta$ we first note that Lemma 2.2 and Lemma 2.6 remain true regardless of the choice of $\delta$. Therefore, the value of $\delta$ affects the behavior of the penalized loss only at the interior points $\boldsymbol{t} \in (0,1)^n$. We would like a choice of $\delta$ such that for all the interior points $\boldsymbol{t} \in (0,1)^n$ the functions $f_\lambda(\boldsymbol{t})$ and $g_\lambda(\boldsymbol{t})$ are well-behaved. When $\delta$ is very small the linear systems that require solving at $\boldsymbol{t} \in (0,1)^n$ may be close to singular and numerical issues can arise more frequently. Moreover, when $\delta$ is large we observe large shifts in the value of the objective approaching a corner point. Our simulations indicate that $\delta = 1$ produces a well-behaved function.

## 4. Numerical Experiments and Results

In this section, we provide numerical examples with real data designed to demonstrate that our proposed continuous optimization method outperforms well-known sampling-based methods for small and large datasets. Moreover, we demonstrate that when it is feasible to run greedy selection, our continuous method exhibits very similar performance.

Numerical experiments were conducted on the small to medium-sized datasets: Residential and Building dataset ($m = 372$, $n = 109$), MNIST1K ($m = 1000$, $n = 784$)[1], Arrhythmia dataset ($m = 452$, $n = 279$), SECOM ($m = 1567$, $n = 591$). Numerical experiments for Nyström land-mark selection were also conducted on the larger datasets: Power Plant dataset ($m = 4$, $n = 9568$), HTRU2 dataset ($m = 8$, $n = 17898$) and Protein dataset ($m = 9$, $n = 45730$). All datasets except MNIST are downloaded from UCI ML Repository (Asuncion & Newman, 2007). All datasets were standardized such that all columns had mean zero and variance equal to one.

For the small to medium-sized datasets, we use the best rank-k approximation factor to compare our method to existing methods (see Figure 2 and Figure 3). The best rank-k approximation factor is given by

$$\text{Approximation Factor} := \frac{\|\mathbf{A} - \widehat{\mathbf{A}}_s\|_F^2}{\|\mathbf{A} - \widehat{\mathbf{G}}\|_F^2}.$$

where $\widehat{\mathbf{A}}_s$ is either the Nyström or CSSP low-rank matrix and $\widehat{\mathbf{G}}$ is the best rank-k approximation computed using the Singular Value Decomposition (SVD) of $\mathbf{A}$.

In these experiments, we compare the proposed continuous landmark selection method executed with SGD ($M = 10$) with the following four well-known methods: Uniform Sam-

---

[1]https://yann.lecun.com/exdb/mnist/

(a) Power Plant  (b) HTRU2  (c) Protein

*Figure 4.* The mean empirical squared Frobenius error $\|\mathbf{K} - \widehat{\mathbf{K}}_{s*}\|_F^2$ over 10 trials for the UCI datasets Power Plant, HTRU2 and Protein for different methods. The kernel matrix $\mathbf{K}$ for all datasets is constructed using the RBF kernel function with $\sigma = 0.5$. Error is plotted on a logarithmic scale.

pling (Williams & Seeger, 2000), Recursive RLS (Ridge Leverage Scores) - Nyström sampling (Musco & Musco, 2017), k-DPP sampling (Derezinski & Mahoney, 2021) and Greedy selection (Farahat et al., 2011; 2013).

For the experiments conducted on the larger datasets (see Figure 4) we exclude the k-DPP sampling and greedy methods as it is either too costly to compute the choice of landmark points or too costly to store the full kernel matrix on a GPU. In our implementation of continuous Nyström landmark selection, we use the KeOps library (Charlier et al., 2021) to efficiently compute MVMs and linear solves on a GPU without ever storing the matrix $\mathbf{K}$, thus negating the need to store any $O(n^2)$ objects. These experiments were run using an NVIDIA Tesla T4 GPU with 16GB memory.

In Figure 2 and Figure 3 we observe the approximation factor for Nyström and CSSP landmark selection with different subset sizes $k$. A lower approximation factor indicates a better approximation and an approximation factor close to one implies near-best-case performance for the given subset size $k$. The results indicate that the continuous optimization method is better than every tested sampling method and is very similar to greedy selection in performance (whenever the greedy selection is feasible). In most cases, for the CSSP, as the proportion of selected columns increases the continuous method starts to marginally outperform the greedy method.

In Figure 4, we observe for all three datasets (Power Plant, HTRU2 and Protein) that the continuous landmark selection achieves better accuracy than the Recursive RLS (Ridge Leverage Scores) - Nyström sampling and Uniform sampling methods. While Recursive RLS sampling (complexity: $O(nk^2)$) and uniform sampling are faster at selecting landmark points, for a fixed $k$ the continuous method obtains a more accurate Nyström approximation. Thus, if a memory budget for the size of the Nyström approximation is given, as is often the case, the continuous method will compute a superior approximation.

## 5. Conclusion

In this paper, we have introduced a novel algorithm that exploits unconstrained continuous optimization to select columns for both the CSSP and Nyström approximation. The algorithm selects columns by minimizing an extended objective which is defined over the hypercube $[0, 1]^n$ rather than iterating over the corner points of the hypercube which correspond to all of the $\binom{n}{k}$ subsets. The extended objective for both the CSSP and Nyström approximation can be minimized via SGD where the gradients are estimated with an unbiased estimator which requires only MVMs with either $\mathbf{X}$ (CSSP) or $\mathbf{K}$ (Nyström). On the real-world examples that we considered in this article, the proposed method has proven to be more accurate without incurring higher computational cost.

## References

Alaoui, A. and Mahoney, M. W. Fast randomized kernel ridge regression with statistical guarantees. *Advances in neural information processing systems*, 28, 2015.

Altschuler, J., Bhaskara, A., Fu, G., Mirrokni, V., Rostamizadeh, A., and Zadimoghaddam, M. Greedy column subset selection: New bounds and distributed algorithms. In *International conference on machine learning*, pp. 2539–2548. PMLR, 2016.

Asuncion, A. and Newman, D. J. Uci machine learning repository, 2007, 2007.

Beale, E., Kendall, M., and Mann, D. The discarding of variables in multivariate analysis. *Biometrika*, 54(3-4): 357–366, 1967.

Bekas, C., Kokiopoulou, E., and Saad, Y. An estimator for

the diagonal of a matrix. *Applied numerical mathematics*, 57(11-12):1214–1229, 2007.

Bien, J., Xu, Y., and Mahoney, M. W. Cur from a sparse optimization viewpoint. *Advances in Neural Information Processing Systems*, 23, 2010.

Calandriello, D., Derezinski, M., and Valko, M. Sampling from a k-dpp without looking at all items. *Advances in Neural Information Processing Systems*, 33:6889–6899, 2020.

Charlier, B., Feydy, J., Glaunes, J. A., Collin, F.-D., and Durif, G. Kernel operations on the gpu, with autodiff, without memory overflows. *J. Mach. Learn. Res.*, 22(74): 1–6, 2021.

Dereziński, M. Fast determinantal point processes via distortion-free intermediate sampling. In *Conference on Learning Theory*, pp. 1029–1049. PMLR, 2019.

Derezinski, M. and Mahoney, M. W. Determinantal point processes in randomized numerical linear algebra. *Notices of the American Mathematical Society*, 68(1):34–45, 2021.

Derezinski, M., Calandriello, D., and Valko, M. Exact sampling of determinantal point processes with sublinear time preprocessing. *Advances in neural information processing systems*, 32, 2019.

Derezinski, M., Khanna, R., and Mahoney, M. W. Improved guarantees and a multiple-descent curve for column subset selection and the nystrom method. *Advances in Neural Information Processing Systems*, 33:4953–4964, 2020.

Deshpande, A. and Vempala, S. Adaptive sampling and fast low-rank matrix approximation. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pp. 292–303. Springer, 2006.

Dietrich, C. R. and Newsam, G. N. Fast and exact simulation of stationary gaussian processes through circulant embedding of the covariance matrix. *SIAM Journal on Scientific Computing*, 18(4):1088–1107, 1997.

Drineas, P., Mahoney, M. W., and Cristianini, N. On the nyström method for approximating a gram matrix for improved kernel-based learning. *journal of machine learning research*, 6(12), 2005.

Farahat, A., Ghodsi, A., and Kamel, M. A novel greedy algorithm for nyström approximation. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pp. 269–277. JMLR Workshop and Conference Proceedings, 2011.

Farahat, A. K., Ghodsi, A., and Kamel, M. S. Efficient greedy feature selection for unsupervised learning. *Knowledge and information systems*, 35(2):285–310, 2013.

Gardner, J., Pleiss, G., Wu, R., Weinberger, K., and Wilson, A. Product kernel interpolation for scalable gaussian processes. In *International Conference on Artificial Intelligence and Statistics*, pp. 1407–1416. PMLR, 2018.

Gittens, A. and Mahoney, M. Revisiting the nystrom method for improved large-scale machine learning. In *International Conference on Machine Learning*, pp. 567–575. PMLR, 2013.

Golub, G. H. and Van Loan, C. F. Matrix computations, 1996.

Hocking, R. R. and Leslie, R. Selection of the best subset in regression analysis. *Technometrics*, 9(4):531–540, 1967.

Hough, J. B., Krishnapur, M., Peres, Y., and Virág, B. Determinantal processes and independence. *Probability surveys*, 3:206–229, 2006.

Hu, R., Sejdinovic, D., and Glaunès, J. A. Giga-scale kernel matrix vector multiplication on gpu. *arXiv preprint arXiv:2202.01085*, 2022.

Li, C., Jegelka, S., and Sra, S. Fast dpp sampling for nystrom with application to kernel methods. In *International Conference on Machine Learning*, pp. 2061–2070. PMLR, 2016.

Martens, J., Sutskever, I., and Swersky, K. Estimating the hessian by back-propagating curvature. In *Proceedings of the 29th International Coference on International Conference on Machine Learning*, pp. 963–970, 2012.

Mathur, A., Moka, S., and Botev, Z. Variance reduction for matrix computations with applications to gaussian processes. In *EAI International Conference on Performance Evaluation Methodologies and Tools*, pp. 243–261. Springer, 2021.

Moka, S., Liquet, B., Zhu, H., and Muller, S. Combss: Best subset selection via continuous optimization. *arXiv preprint arXiv:2205.02617*, 2022.

Musco, C. and Musco, C. Recursive sampling for the nystrom method. *Advances in neural information processing systems*, 30, 2017.

Robbins, H. and Monro, S. A stochastic approximation method. *The annals of mathematical statistics*, pp. 400–407, 1951.

Ryan, J. P., Ament, S. E., Gomes, C. P., and Damle, A. The fast kernel transform. In *International Conference on Artificial Intelligence and Statistics*, pp. 11669–11690. PMLR, 2022.

Shitov, Y. Column subset selection is np-complete. *Linear Algebra and its Applications*, 610:52–58, 2021.

Tibshirani, R. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.

Williams, C. and Seeger, M. Using the nyström method to speed up kernel machines. *Advances in neural information processing systems*, 13, 2000.

Yuan, M. and Lin, Y. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68 (1):49–67, 2006.

## A. Proofs

*Proof of Lemma 2.2.* The following proof follows similar arguments to that of *Theorem 1* in (Moka et al., 2022). Given that the pseudo-inverse of a matrix after a permutation of rows (respectively, columns) is identical to the matrix obtained by applying the same permutation on columns (respectively, rows) on the pseudo-inverse, we assume without loss of generality that all the zero-elements $s \in \{0, 1\}^n$ appear at the end, in the form,

$$s = (s_1, \ldots s_l, 0, \ldots, 0).$$

where $l$ is equal to the number of non-zeros in $s \in \{0, 1\}^n$. Then, $\widetilde{\mathbf{P}}(s)$ is given by the block-wise matrix,

$$\widetilde{\mathbf{P}}(s) = \begin{bmatrix} \mathbf{X}_{[s]} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{X}_{[s]}^\top \mathbf{X}_{[s]} & \mathbf{0} \\ \mathbf{0} & \delta \mathbf{I} \end{bmatrix}^\dagger \begin{bmatrix} \mathbf{X}_{[s]}^\top \\ \mathbf{0} \end{bmatrix}. \tag{4}$$

It is easy to verify, when the matrices $\mathbf{A}_1$ and $\mathbf{A}_2$ are square, the block-diagonal pseudo-inverse $\begin{bmatrix} \mathbf{A}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_2 \end{bmatrix}^\dagger = \begin{bmatrix} \mathbf{A}_1^\dagger & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_2^\dagger \end{bmatrix}$.

Therefore (4) reduces to,

$$\widetilde{\mathbf{P}}(s) = \begin{bmatrix} \mathbf{X}_{[s]} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{X}_{[s]}^\top \mathbf{X}_{[s]}^\dagger & \mathbf{0} \\ \mathbf{0} & \delta^{-1}\mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{X}_{[s]}^\top \\ \mathbf{0} \end{bmatrix}$$

$$= \mathbf{X}_{[s]} \left( \mathbf{X}_{[s]}^\top \mathbf{X}_{[s]} \right)^\dagger \mathbf{X}_{[s]}^\top$$

$$= \mathbf{X}_{[s]} \mathbf{X}_{[s]}^\dagger.$$

□

*Proof of Lemma 2.2.* To obtain the gradient $f_\lambda(\mathbf{t})$ for $\mathbf{t} \in (0, 1)^n$ we first simplify the term $-\operatorname{tr}\left[\mathbf{X}^\top \widetilde{\mathbf{P}}(\mathbf{t}) \mathbf{X}\right]$ by letting $\mathbf{K} = \mathbf{X}^\top \mathbf{X}$, $\mathbf{Z} = \mathbf{K} - \delta \mathbf{I}$ and $\mathbf{L}_t = \mathbf{TZT} + \delta \mathbf{I}$. Then, we have,

$$-\operatorname{tr}\left[\mathbf{X}^\top \widetilde{\mathbf{P}}(\mathbf{t}) \mathbf{X}\right] = -\operatorname{tr}\left[\mathbf{KTL}_t^{-1}\mathbf{TK}\right].$$

Using matrix calculus, for any $j = 1, \ldots, n$, we have the partial derivative,

$$\frac{\partial}{\partial t_j}\left[-\operatorname{tr}\left(\mathbf{X}^\top \widetilde{\mathbf{P}}(\mathbf{t}) \mathbf{X}\right)\right] = -\operatorname{tr}\left(\mathbf{K}\frac{\partial\left[\mathbf{TL}_t^{-1}\mathbf{T}\right]}{\partial t_j}\mathbf{K}\right). \tag{5}$$

Let $\mathbf{E}_j$ be the square matrix of dimension $n \times n$ with 1 at position $(j, j)$ and 0 everywhere else. Then $\frac{\partial \mathbf{T}}{\partial t_j} = \mathbf{E}_j$ and we have,

$$\frac{\partial\left[\mathbf{TL}_t^{-1}\mathbf{T}\right]}{\partial t_j} = \mathbf{E}_j\mathbf{L}_t^{-1}\mathbf{T} + \mathbf{T}\frac{\partial\left[\mathbf{L}_t^{-1}\mathbf{T}\right]}{\partial t_j}. \tag{6}$$

Furthermore,

$$\frac{\partial\left[\mathbf{L}_t^{-1}\mathbf{T}\right]}{\partial t_j} = \frac{\partial\left[\mathbf{L}_t^{-1}\right]}{\partial t_j}\mathbf{T} + \mathbf{L}_t^{-1}\mathbf{E}_j, \tag{7}$$

and using the derivative of an invertible matrix we have,

$$\frac{\partial\left[\mathbf{L}_t^{-1}\right]}{\partial t_j} = -\mathbf{L}_t^{-1}\frac{\partial \mathbf{L}_t}{\partial t_j}\mathbf{L}_t^{-1}, \tag{8}$$

and,

$$\frac{\partial \mathbf{L}_t}{\partial t_j} = \mathbf{E}_j\mathbf{ZT} + \mathbf{TZE}_j. \tag{9}$$

Substituting (9) → (8) → (7) → (6) → (5) we obtain the expression,

$$\frac{\partial}{\partial t_j}\left[-\operatorname{tr}\left(\mathbf{X}^\top \widetilde{\mathbf{P}}(\mathbf{t}) \mathbf{X}\right)\right] = -\operatorname{tr}\left[\mathbf{KE}_j\mathbf{L}_t^{-1}\mathbf{TK} + \mathbf{KTL}_t^{-1}\mathbf{E}_j\mathbf{K} - \mathbf{KTL}_t^{-1}\mathbf{E}_j\mathbf{ZTL}_t^{-1}\mathbf{TK} - \mathbf{KTL}_t^{-1}\mathbf{TZE}_j\mathbf{L}_t^{-1}\mathbf{TK}\right]. \tag{10}$$

In order to simplify this expression we consider the following fact. If we have the matrices $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{n \times n}$ and $\mathbf{B} = (b_{ij}) \in \mathbb{R}^{n \times n}$ then, $\operatorname{tr}(\mathbf{A} \mathbf{E}_j \mathbf{B}) = (\mathbf{B} \mathbf{A})_{jj} = \sum_{i=1}^{n} a_{ij} b_{ji}$. Using this, we obtain,

$$\frac{\partial}{\partial t_j} \left[ - \operatorname{tr}\left( \mathbf{X}^\top \widetilde{\mathbf{P}}(t) \mathbf{X} \right) \right] = - \left[ \mathbf{L}_t^{-1} \mathbf{T} \mathbf{K}^2 \right]_{jj} - \left[ \mathbf{K}^2 \mathbf{T} \mathbf{L}_t^{-1} \right]_{jj} + \left[ \mathbf{Z} \mathbf{T} \mathbf{L}_t^{-1} \mathbf{T} \mathbf{K}^2 \mathbf{T} \mathbf{L}_t^{-1} \right]_{jj} + \left[ \mathbf{L}_t^{-1} \mathbf{T} \mathbf{K}^2 \mathbf{T} \mathbf{L}_t^{-1} \mathbf{Z} \right]_{jj}$$

$$= 2 \left[ \mathbf{L}_t^{-1} \mathbf{T} \mathbf{K}^2 \left( \mathbf{T} \mathbf{L}_t^{-1} \mathbf{Z} - \mathbf{I} \right) \right]_{jj},$$

since the matrices $\mathbf{L}_t^{-1}$, $\mathbf{K}$, $\mathbf{Z}$ and $\mathbf{T}$ are all symmetric. Considering the partial derivative of the penalty term is $\frac{\partial}{\partial t_j} \left[ \lambda \sum_i t_i \right] = \lambda$ we have the following expression for the gradient vector of $f_\lambda(t)$,

$$\nabla f_\lambda(t) = 2 \operatorname{Diag} \left[ \mathbf{L}_t^{-1} \mathbf{T} \mathbf{K}^2 \left( \mathbf{T} \mathbf{L}_t^{-1} \mathbf{T} \mathbf{Z} - \mathbf{I} \right) \right] + \lambda \mathbf{1}.$$

☐

*Proof of Lemma 2.6.* For reasons outlined in the proof of Lemma 2.2 we assume without loss of generality that all the zero-elements in $s \in \{0, 1\}^n$ appear at the end, in the form,

$$s = (s_1, \ldots s_l, 0, \ldots, 0).$$

where $l$ is equal to the number of non-zeros in $s \in \{0, 1\}^n$. Then, $\widetilde{K}(s)$ is given by the block-wise matrix,

$$\widetilde{\mathbf{P}}(s) = \begin{bmatrix} \mathbf{K}_{[s]} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{K}_{[s,s]} & \mathbf{0} \\ \mathbf{0} & \delta \mathbf{I} \end{bmatrix}^\dagger \begin{bmatrix} \mathbf{K}_{[s]}^\top \\ \mathbf{0} \end{bmatrix}. \tag{11}$$

Using the block-diagonal pseudo-inverse formula we have,

$$\widetilde{\mathbf{P}}(s) = \begin{bmatrix} \mathbf{K}_{[s]} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{K}_{[s,s]}^\dagger & \mathbf{0} \\ \mathbf{0} & \delta^{-1} \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{K}_{[s]}^\top \\ \mathbf{0} \end{bmatrix}$$

$$= \mathbf{K}_{[s]} \mathbf{K}_{[s,s]}^\dagger \mathbf{K}_{[s]}^\top.$$

☐

*Proof of Lemma 2.7.* For any positive semi-definite kernel matrix $\mathbf{K} \in \mathbb{R}^{n \times n}$, a decomposition of the form $\mathbf{K} = \mathbf{X}^\top \mathbf{X}$ where $\mathbf{X} \in \mathbb{R}^{n \times n}$ always exists. Therefore, the function $\widetilde{K}(t)$ can be written as,

$$\widetilde{K}(t) = \mathbf{X}^\top \mathbf{X} \mathbf{T} \left[ \mathbf{T} \mathbf{X}^\top \mathbf{X} \mathbf{T} + \delta (\mathbf{I} - \mathbf{T}^2) \right]^\dagger \mathbf{T} \mathbf{X}^\top \mathbf{X}$$

$$= \mathbf{X}^\top \widetilde{\mathbf{P}}(t) \mathbf{X}^\top.$$

From Lemma 2.3 we know that the function $\widetilde{\mathbf{P}}(t)$ is continuous over $[0, 1]^n$. Since $\mathbf{X}$ is not a function of $t$ we conclude that $\widetilde{K}(t)$ is also continuous over $[0, 1]^n$. ☐

*Proof of Lemma 2.8.* To obtain the gradient $g_\lambda(t)$ for $t \in (0, 1)^n$ we first simplify the term $\|\widetilde{K}(t) - \mathbf{K}\|_F^2$. Since $\widetilde{K}(t)$ and $\mathbf{K}$ are symmetric, we have the expansion,

$$\|\widetilde{K}(t) - \mathbf{K}\|_F^2 = \operatorname{tr}\left[ \left( \widetilde{K}(t) - \mathbf{K} \right)^2 \right]$$

$$= \operatorname{tr}\left[ \left( \widetilde{K}(t) \right)^2 - 2 \widetilde{K}(t) \mathbf{K} + \mathbf{K}^2 \right].$$

Therefore,

$$\frac{\partial}{\partial t_j}\left[\|\widetilde{\mathbf{K}}(\boldsymbol{t}) - \mathbf{K}\|_F^2\right] = \operatorname{tr}\left(\frac{\partial}{\partial t_j}\left[\left(\widetilde{\mathbf{K}}(\boldsymbol{t})\right)^2\right]\right) - \operatorname{tr}\left(\frac{\partial}{\partial t_j}\left[2\widetilde{\mathbf{K}}(\boldsymbol{t})\mathbf{K}\right]\right)$$

$$= 2\operatorname{tr}\left(\frac{\partial\widetilde{\mathbf{K}}(\boldsymbol{t})}{\partial t_j}\widetilde{\mathbf{K}}(\boldsymbol{t})\right) - 2\operatorname{tr}\left(\frac{\partial\widetilde{\mathbf{K}}(\boldsymbol{t})}{\partial t_j}\mathbf{K}\right)$$

$$= 2\operatorname{tr}\left(\frac{\partial\widetilde{\mathbf{K}}(\boldsymbol{t})}{\partial t_j}\mathbf{D}\right)$$

where $\mathbf{D} = \widetilde{\mathbf{K}}(\boldsymbol{t}) - \mathbf{K}$. Factorize $\mathbf{K}$ as $\mathbf{K} = \mathbf{X}^\top\mathbf{X}$ and let $\mathbf{Z} = \mathbf{K} - \delta\mathbf{I}$ and $\mathbf{L_t} = \mathbf{TZT} + \delta\mathbf{I}$. Then, notice that the derivative $\frac{\partial\widetilde{\mathbf{K}}(\boldsymbol{t})}{\partial t_j} = \frac{\partial[\mathbf{X}^\top\widetilde{\mathbf{P}}(\boldsymbol{t})\mathbf{X}]}{\partial t_j}$ is the same expression that we derived in the proof of Lemma 2.4, see (10). Substituting this expression in for $\frac{\partial\widetilde{\mathbf{K}}(\boldsymbol{t})}{\partial t_j}$, we obtain,

$$\frac{\partial}{\partial t_j}\left[\|\widetilde{\mathbf{K}}(\boldsymbol{t}) - \mathbf{K}\|_F^2\right] = 2\operatorname{tr}\left[\mathbf{K}\mathbf{E}_j\mathbf{L_t}^{-1}\mathbf{TKD} + \mathbf{KTL_t}^{-1}\mathbf{E}_j\mathbf{KD} - \mathbf{KTL_t}^{-1}\mathbf{E}_j\mathbf{ZTL_t}^{-1}\mathbf{TKD} - \mathbf{KTL_t}^{-1}\mathbf{TZE}_j\mathbf{L_t}^{-1}\mathbf{TKD}\right].$$

Once again, using the fact that $\operatorname{tr}(\mathbf{AE}_j\mathbf{B}) = (\mathbf{BA})_{jj} = \sum_{i=1}^n a_{ij}b_{ji}$ for matrices $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{n\times n}$ and $\mathbf{B} = (b_{ij}) \in \mathbb{R}^{n\times n}$, we obtain,

$$\frac{\partial}{\partial t_j}\left[\|\widetilde{\mathbf{K}}(\boldsymbol{t}) - \mathbf{K}\|_F^2\right] = 2\left[\left(\mathbf{L_t}^{-1}\mathbf{TKDK}\right)_{jj} + \left(\mathbf{KDKTL_t}^{-1}\right)_{jj} - \left(\mathbf{ZTL_t}^{-1}\mathbf{TKDKTL_t}^{-1}\right)_{jj} - \left(\mathbf{L_t}^{-1}\mathbf{TKDKTL_t}^{-1}\mathbf{TZ}\right)_{jj}\right]$$

(12)

$$= 4\left[\left(\mathbf{L_t}^{-1}\mathbf{TKDK}\right)_{jj} - \left(\mathbf{L_t}^{-1}\mathbf{TKDKTL_t}^{-1}\mathbf{TZ}\right)_{jj}\right],$$ 

(13)

since the matrices $\mathbf{L_t}^{-1}$, $\mathbf{K}$, $\mathbf{Z}$ and $\mathbf{T}$ are all symmetric. Considering the partial derivative of the penalty term is $\frac{\partial}{\partial t_j}\left[\lambda\sum_i t_i\right] = \lambda$ we have the following expression for the gradient vector of $g_\lambda(\boldsymbol{t})$,

$$\nabla g_\lambda(\boldsymbol{t}) = 4\operatorname{Diag}\left[\mathbf{L_t}^{-1}\mathbf{TKDK}\left(\mathbf{I} - \mathbf{TL_t}^{-1}\mathbf{TZ}\right)\right] + \lambda\mathbf{1}.$$

□ *Proof of Lemma 3.1.* From Lemma 2.4, we have,

$$\nabla f_\lambda(\boldsymbol{t}) = 2\operatorname{Diag}\left[\mathbf{L_t}^{-1}\mathbf{TK}^2\left(\mathbf{TL_t}^{-1}\mathbf{TZ} - \mathbf{I}\right)\right] + \lambda\mathbf{1}$$
$$= 2\left[\operatorname{Diag}\left(\mathbf{L_t}^{-1}\mathbf{TK}^2\mathbf{TL_t}^{-1}\mathbf{TZ}\right) - \operatorname{Diag}\left(\mathbf{L_t}^{-1}\mathbf{TK}^2\right)\right] + \lambda\mathbf{1}$$
$$= 2\boldsymbol{\alpha} + \lambda\mathbf{1},$$

where $\boldsymbol{\alpha} = \operatorname{Diag}\left(\mathbf{L_t}^{-1}\mathbf{TK}^2\mathbf{TL_t}^{-1}\mathbf{TZ}\right) - \operatorname{Diag}\left(\mathbf{L_t}^{-1}\mathbf{TK}^2\right)$. To obtain an unbiased estimator for $\boldsymbol{\alpha}$, we use the factorized estimator $\hat{\ell}$ for the diagonal of a square matrix. Recall, to estimate the diagonal of the matrix $\mathbf{A} = \mathbf{BC}^\top$ where $\mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathbb{R}^{n\times n}$ we let $\boldsymbol{z} \in \mathbb{R}^n$ be a random vector sampled from the Rademacher distribution. Then the unbiased estimate for $\operatorname{Diag}(\mathbf{A})$ is $\hat{\ell} = \mathbf{B}\boldsymbol{z} \odot \mathbf{C}\boldsymbol{z}$ (see (Martens et al., 2012; Mathur et al., 2021) for proof and analysis). We factorize the matrices in $\boldsymbol{\alpha}$ so that,

$$\boldsymbol{\alpha} = \operatorname{Diag}\left(\overbrace{\mathbf{L_t}^{-1}\mathbf{TK}}^{\mathbf{B}_1}\overbrace{\mathbf{KTL_t}^{-1}\mathbf{TZ}}^{\mathbf{C}_1^\top}\right) - \operatorname{Diag}\left(\overbrace{\mathbf{L_t}^{-1}\mathbf{TK}}^{\mathbf{B}_2}\overbrace{\mathbf{K}}^{\mathbf{C}_2^\top}\right)$$

Then, the factorized estimator for $\boldsymbol{\alpha}$ is given by,

$$\boldsymbol{\phi} = \mathbf{L_t}^{-1}\mathbf{TK}\boldsymbol{z} \odot \mathbf{ZTL_t}^{-1}\mathbf{TK}\boldsymbol{z} - \mathbf{L_t}^{-1}\mathbf{TK}\boldsymbol{z} \odot \mathbf{K}\boldsymbol{z},$$

where $\mathbb{E}[\boldsymbol{\phi}] = \boldsymbol{\alpha}$ and $\boldsymbol{z} \in \mathbb{R}^n$ is a Rademacher random variable. If we compute the following variables: $(1)\,\boldsymbol{a} = \mathbf{K}\boldsymbol{z}$, $(2)\,\boldsymbol{b} = \mathbf{L_t}^{-1}(\boldsymbol{t} \odot \boldsymbol{a})$, we have,

$$(1)\,\boldsymbol{a} = \mathbf{K}\boldsymbol{z}, \quad (2)\,\boldsymbol{b} = \mathbf{L_t}^{-1}\mathbf{TK}\boldsymbol{z}$$

and $\phi$ simplifies to,

$$\phi = \boldsymbol{b} \odot \mathbf{Z}(\boldsymbol{t} \odot \boldsymbol{b}) - \boldsymbol{a} \odot \boldsymbol{b}.$$

Therefore, for $\boldsymbol{t} \in (0,1)^n$ we have,

$$\nabla f_\lambda(\boldsymbol{t}) = 2\mathbb{E}\left[\phi\right] + \lambda \mathbf{1}.$$

$\square$

*Proof of Lemma 3.2.* From Lemma 2.8 and Equation (12), we have,

$$
\begin{aligned}
\nabla g_\lambda(\boldsymbol{t}) &= 4\,\mathrm{Diag}\left[\mathbf{L}_t^{-1}\mathbf{TKDK}\left(\mathbf{I} - \mathbf{TL}_t^{-1}\mathbf{TZ}\right)\right] + \lambda\mathbf{1}. \\
&= 2\left[\mathrm{Diag}\left(\mathbf{L}_t^{-1}\mathbf{TKDK}\right) + \mathrm{Diag}\left(\mathbf{KDKTL}_t^{-1}\right) - \mathrm{Diag}\left(\mathbf{ZTL}_t^{-1}\mathbf{TKDKTL}_t^{-1}\right) - \mathrm{Diag}\left(\mathbf{L}_t^{-1}\mathbf{TKDKTL}_t^{-1}\mathbf{TZ}\right)\right] + \lambda\mathbf{1} \\
&= 2\boldsymbol{\beta} + \lambda\mathbf{1}.
\end{aligned}
$$

To obtain an unbiased estimator for $\boldsymbol{\beta}$, we once again use the factorized estimator for the diagonal of a square matrix. We factorize the matrices in $\boldsymbol{\beta}$ so that,

$$\boldsymbol{\beta} = \mathrm{Diag}\left(\overbrace{\mathbf{L}_t^{-1}\mathbf{TKD}}^{\mathbf{B}_1}\overbrace{\mathbf{K}}^{\mathbf{C}_1^\top}\right) + \mathrm{Diag}\left(\overbrace{\mathbf{KD}}^{\mathbf{B}_2}\overbrace{\mathbf{KTL}_t^{-1}}^{\mathbf{C}_2^\top}\right) - \mathrm{Diag}\left(\overbrace{\mathbf{ZTL}_t^{-1}\mathbf{TKD}}^{\mathbf{B}_3}\overbrace{\mathbf{KTL}_t^{-1}}^{\mathbf{C}_3^\top}\right) - \mathrm{Diag}\left(\overbrace{\mathbf{L}_t^{-1}\mathbf{TKD}}^{\mathbf{B}_4}\overbrace{\mathbf{KTL}_t^{-1}\mathbf{TZ}}^{\mathbf{C}_4}\right).$$

Then, the factorized estimator for $\boldsymbol{\beta}$ is given by,

$$\psi = \mathbf{L}_t^{-1}\mathbf{TKD}\boldsymbol{z} \odot \mathbf{K}\boldsymbol{z} + \mathbf{KD}\boldsymbol{z} \odot \mathbf{L}_t^{-1}\mathbf{TK}\boldsymbol{z} - \mathbf{ZTL}_t^{-1}\mathbf{TKD}\boldsymbol{z} \odot \mathbf{L}_t^{-1}\mathbf{TK}\boldsymbol{z} - \mathbf{L}_t^{-1}\mathbf{TKD}\boldsymbol{z} \odot \mathbf{ZTL}_t^{-1}\mathbf{TK}\boldsymbol{z}.$$

where $\mathbb{E}[\psi] = \boldsymbol{\beta}$ and $\boldsymbol{z} \in \mathbb{R}^n$ is a Rademacher random variable. Recall that $\mathbf{D} = \widetilde{\mathbf{K}}(\boldsymbol{t}) - \mathbf{K}$ and $\widetilde{\mathbf{K}}(\boldsymbol{t}) = \mathbf{KTL}_t^{-1}\mathbf{TK}$. Then, if we compute the following variables: (1) $\boldsymbol{a} = \mathbf{K}\boldsymbol{z}$, (2) $\boldsymbol{b} = \mathbf{L}_t^{-1}(\boldsymbol{t} \odot \boldsymbol{a})$, (3) $\boldsymbol{c} = \mathbf{K}(\boldsymbol{t} \odot \boldsymbol{b}) - \boldsymbol{a}$, (4) $\boldsymbol{d} = \mathbf{K}\boldsymbol{c}$, (5) $\boldsymbol{e} = \mathbf{L}_t^{-1}(\boldsymbol{t} \odot \boldsymbol{d})$, we have,

$$(1)\,\boldsymbol{a} = \mathbf{K}\boldsymbol{z}, \quad (2)\,\boldsymbol{b} = \mathbf{L}_t^{-1}\mathbf{TK}\boldsymbol{z}, \quad (3)\,\boldsymbol{c} = \mathbf{D}\boldsymbol{z}, \quad (4)\,\boldsymbol{d} = \mathbf{KD}\boldsymbol{z} \quad (5)\,\boldsymbol{e} = \mathbf{L}_t^{-1}\mathbf{TKD}\boldsymbol{z},$$

and $\psi$ simplifies to,

$$\psi = \boldsymbol{b} \odot \boldsymbol{d} + \boldsymbol{a} \odot \boldsymbol{e} - \boldsymbol{e} \odot \mathbf{Z}(\boldsymbol{t} \odot \boldsymbol{b}) - \boldsymbol{b} \odot \mathbf{Z}(\boldsymbol{t} \odot \boldsymbol{e}).$$

Hence, we have the expression for the gradient,

$$\nabla g_\lambda(\boldsymbol{t}) = 2\mathbb{E}\left[\psi\right] + \lambda\mathbf{1}.$$

$\square$ *Proof of Lemma 2.3* For the same reasons outlined in the proof of Lemma 2.2 we assume without loss of generality that all the zero-elements in $\boldsymbol{t} \in \{0,1\}^n$ appear at the end, in the form,

$$\boldsymbol{t} = (t_1, \ldots\, t_l, 0, \ldots, 0).$$

Then $\mathbf{L}_t$ is given by,

$$\mathbf{L}_t = \begin{bmatrix} (\mathbf{T})_+(\mathbf{K})_+(\mathbf{T})_+ + \delta(\mathbf{I} - (\mathbf{T})_+^2) & \mathbf{0} \\ \mathbf{0} & \delta\mathbf{I} \end{bmatrix} = \begin{bmatrix} (\mathbf{L}_t)_+ & \mathbf{0} \\ \mathbf{0} & \delta\mathbf{I} \end{bmatrix}$$

and since $(\boldsymbol{t})_+ \in [0,1)^l$ the matrix $(\mathbf{L}_t)_+$ is invertible. Therefore,

$$
\begin{aligned}
\mathbf{L}_t^{-1}(\boldsymbol{t} \odot \boldsymbol{r}) &= \begin{bmatrix} (\mathbf{L}_t)_+^{-1} & \mathbf{0} \\ \mathbf{0} & \delta^{-1}\mathbf{I} \end{bmatrix}\begin{bmatrix} (\boldsymbol{t})_+ \odot (\boldsymbol{r})_+ \\ \mathbf{0} \end{bmatrix} \\
&= \begin{bmatrix} (\mathbf{L}_t)_+^{-1}\left((\boldsymbol{t})_+ \odot (\boldsymbol{r})_+\right) \\ \mathbf{0} \end{bmatrix}.
\end{aligned}
$$

$\square$