

# Sparse Contextual CDF Regression

Anonymous authors

Paper under double-blind review

## Abstract

Estimating cumulative distribution functions (CDFs) of context dependent random variables is a central statistical task underpinning numerous applications in machine learning and economics. In this work, we extend a recent line of theoretical inquiry into this domain by analyzing the problem of *sparse contextual CDF regression*, wherein data points are sampled from a convex combination of  $s$  context dependent CDFs chosen from a set of  $d$  basis functions. We show that adaptations of several canonical regression methods serve as tractable estimators in this functional sparse regression setting under standard assumptions on the conditioning of the basis functions. In particular, given  $n$  data samples, we prove estimation error upper bounds of  $\tilde{O}(\sqrt{s/n})$  for functional versions of the lasso and Dantzig selector estimators, and  $\tilde{O}(\sqrt{s}/\sqrt[4]{n})$  for a functional version of the elastic net estimator. Our results match the corresponding error bounds for finite dimensional regression and improve upon CDF ridge regression which has  $\tilde{O}(\sqrt{d/n})$  sample complexity. Finally, we obtain a matching information-theoretic lower bound which establishes the minimax optimality of the lasso and Dantzig selector estimators up to logarithmic factors.

## 1 Introduction

The estimation of cumulative distribution functions (CDFs) is a classical problem in mathematical statistics stemming back to the Glivenko-Cantelli theorem (Cantelli, 1933; Glivenko, 1933; Devroye et al., 2013), which states that empirical CDFs constructed from independent samples of a single random variable converge uniformly to the random variable’s true CDF. Subsequent classical research in this area has focused on deriving tight non-asymptotic sample complexity results in terms of the Kolmogorov-Smirnov distance among others, such as the Dvoretzky-Kiefer-Wolfowitz inequality (Dvoretzky et al., 1956) and improved bounds by Massart (1990).

Motivated by applications to modern learning tasks such as contextual bandits and Markov decision processes (Huang et al., 2021; 2022), a recent line of research (Zhang et al., 2024) introduced the problem of *context dependent* CDF estimation, which requires a learner to simultaneously estimate a (possibly infinite) family of CDFs parameterized by some context variable. As an initial simplification, the authors considered the restricted setting of *contextual CDF regression*, wherein the true contextual CDF is a convex combination of  $d$  context-dependent basis functions. The authors generalized the classical ridge regression method (Hoerl & Kennard, 1970; Abbasi-Yadkori et al., 2011) to this functional regression problem and derived a tight  $\tilde{O}(\sqrt{d/n})$  estimation error bound given  $n$  samples in a variety of data generation settings. However, when  $d$  is large, specifically in unstructured settings where a massive set of potential CDFs are considered, ridge regression utilizes all CDF basis functions for the purpose of estimation without regard to their relevance.

In this paper, as a further step towards developing general algorithms for contextual CDF estimation, we propose *sparse regression* and *basis selection* techniques for the aforementioned CDF regression problem based on functional versions of lasso (Tibshirani, 1996), elastic net (Zou & Hastie, 2005), and Dantzig selector (Candes & Tao, 2007) methods. Crucially, all of our techniques achieve estimation bounds with no polynomial dependence on  $d$ , allowing accurate recovery of the true contextual CDF from bases containing exponentially many irrelevant functions. We also establish minimax optimality results for sparse CDF regression.

## 2 Outline

We briefly delineate the structure of our paper. In Section 3, we outline some applications of CDFs in machine learning and economics, and provide an overview of related functional regression schemes in the previous scientific literature. We introduce the sparse contextual CDF regression problem under present investigation in Section 4 and define relevant notation in Appendix A. In Section 5, we state our main contributions and discuss the overarching proof techniques used throughout the paper. Section 6 provides the formal derivation of our upper bound on lasso estimation error in the fixed design setting. We include numerical simulations of our data generation and parameter estimation processes in Section 7. We defer remaining proofs and technical details to Appendices B to E, and provide examples of CDF bases which satisfy the preconditions for our main results in Appendix F.

## 3 Related Work

Broadly speaking, CDFs underpin the computation of risk functionals which inform decision-making in mathematical finance and actuarial science. For example, generic law invariant risk functions such as conditional value-at-risk (Rockafellar et al., 2000; Artzner et al., 1999) are parameterized by CDFs, and a notion of distortion risk measure (Wirch & Hardy, 2001) arises from the composition of a CDF with a distortion function. Coherent risk measures are instrumental in portfolio management and optimization (Krokhmal, 2007) and can be formulated in terms of CDFs (Shapiro et al., 2014). Lastly, spectral risk measures are computed as weighted averages of outcomes (Acerbi, 2002) and hence depend on the entire trajectory of the underlying random variable’s CDF.

Similar considerations exist in further scopes of learning theory, where the aforementioned risk functionals are incorporated into supervised learning tasks (Liu et al., 2022) and multi-armed bandit problems (Cassel et al., 2023) to model fairness, risk aversion, and distribution shift (Wong et al., 2022). The recently proposed off-policy risk assessment framework (Huang et al., 2021) includes CDF estimation as a key building block and has been applied to contextual bandit problems and Markov decision processes (Huang et al., 2022). Other work in this vein has employed the mean-variance (Sani et al., 2013; Zimin et al., 2014) and value-at-risk (Vakili & Zhao, 2015) paradigms in the multi-armed bandit setting to analyze risk-reward trade-offs. Furthermore, cumulative prospect theory is intimately tied to risk distortion and the ensuing dependence on CDF estimation (Prashanth et al., 2016), and has found relevance in reinforcement learning and stochastic optimization (Jie et al., 2018).

We touch on the well-established precedent in the scientific community for framing function estimation through the lens of linear regression and basis selection. Romo et al. (2013) generalized lasso variable selection to linear models with scalar regressors and functional responses to develop interpretable analyses of car accident data. A related method exploiting feature sparsity and output function smoothness was proposed in Barber et al. (2017) to perform genome-wide association studies. Linear models with functional regressors and scalar responses have been used in genomics, MRI data analysis, and chemometrics, spurring the development of functional group-lasso (Pannu & Billor, 2017) and wavelet-based lasso methods (Zhao et al., 2012; 2015) for such models. Prior studies on genetic regulatory networks (Hong & Lian, 2011) have incorporated function-on-function regression models with  $\ell^1$ -regularization to encode sparsity in pairwise interactions between genes. Lastly, recent developments on lasso estimation for function-on-function regression (Centofanti et al., 2022; Maranzano et al., 2023) find downstream application in geostatistical models and mortality data analysis.

Compared to the prior literature, the main novelty of our contributions lies in the simultaneous estimation of an entire family of CDFs. In this sense, our work may be interpreted as a generalization of the canonical mixture model with known basis distributions (Murphy, 2012, Section 11.2) to the context dependent setting. For maximum generality, we focus on CDF estimation instead of PDFs or quantiles which require restrictions for existence and well-definedness. We formulate estimation through the lens of functional linear regression and show that our task reduces to integral computation and finite dimensional optimization, making our approach computationally efficient using existing numerical methods. Lastly, our derivations utilize intrinsic

properties of CDFs instead of performing discretization to avoid introducing approximation error, in contrast to the approach taken in [Hong & Lian \(2011\)](#) among others.

## 4 Model and Setup

In this section, we formally define the sparse contextual CDF regression problem under investigation in this paper. (Recall that we utilize the notation defined in [Appendix A](#).) Let  $\mathcal{X}$  denote a general context space. We refer to a function  $f(x, t) : \mathcal{X} \times \mathbb{R} \rightarrow [0, 1]$  as a *contextual CDF* if for any  $x \in \mathcal{X}$ ,  $f(x, \cdot)$  is a valid CDF for a real-valued random variable with range contained in some set  $S \subseteq \mathbb{R}$ . Let  $\mathbf{m}$  be any probability measure on  $S$ . Let  $\{\phi_1, \dots, \phi_d\}$  denote a fixed basis of  $d \in \mathbb{N}$  contextual CDFs indexed by  $i \in [d]$ . For convenience, we often conceptualize this basis as a single vector-valued function  $\Phi : \mathcal{X} \times \mathbb{R} \rightarrow [0, 1]^d$  defined by  $[\Phi(x, t)]_i = \phi_i(x, t)$  for all  $i \in [d]$ .

Let  $F : \mathcal{X} \times \mathbb{R} \rightarrow [0, 1]$  denote the true contextual CDF we aim to recover. Following the precedent established in [Zhang et al. \(2024\)](#), we assume that  $F$  is a convex combination of the basis functions  $\{\phi_1, \dots, \phi_d\}$ . Hence, the problem reduces to recovering the true parameter vector  $\theta_* \in \Delta^{d-1}$  such that

$$\forall x \in \mathcal{X}, \forall t \in \mathbb{R}, F(x, t) = \theta_*^\top \Phi(x, t).$$

Let  $S_* = \text{supp}(\theta_*)$  denote the support of the true parameter. For notational simplicity, let  $s = |S_*| = \|\theta_*\|_0$ . Let  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$  be a set of  $n \in \mathbb{N}$  observed samples indexed by  $j \in [n]$ , generated according to the *fixed design* or *random design* settings described below:

- **Fixed design.** For each  $j \in [n]$ , the context variable  $x^{(j)} \in \mathcal{X}$  is fixed a priori, and the response variable  $y^{(j)} \in \mathbb{R}$  is independently sampled from the CDF  $F(x^{(j)}, \cdot)$ .
- **Random design.** For each  $j \in [n]$ , the context variable  $x^{(j)} \in \mathcal{X}$  is independently sampled from an unknown probability distribution  $P_X^{(j)}$  over  $\mathcal{X}$ . Then, conditional on  $x^{(j)}$ , the response variable  $y^{(j)} \in \mathbb{R}$  is independently sampled from the CDF  $F(x^{(j)}, \cdot)$ .

These design settings extend the data generation processes described in [Abbasi-Yadkori et al. \(2011\)](#) and [Hsu et al. \(2012\)](#). Let  $\Phi_j(t) = \Phi(x^{(j)}, t)$  denote the CDF basis at the  $j$ th context variable, and let  $x^{1:n} = (x^{(1)}, \dots, x^{(n)})$  denote a collection of sampled variables, with analogous definitions for  $y^{1:n}$  and  $(x, y)^{1:n}$ . Let  $U_n \in \mathbb{R}^{d \times d}$  denote the *empirical  $n$ -sample Gramian matrix* given by

$$\forall i, i' \in [d], [U_n]_{i, i'} = \frac{1}{n} \sum_{j=1}^n \langle [\Phi_j]_i, [\Phi_j]_{i'} \rangle \quad \therefore U_n = \frac{1}{n} \sum_{j=1}^n \int_S \Phi_j \Phi_j^\top d\mathbf{m}.$$

In the random design setting, let  $\Sigma_n \in \mathbb{R}^{d \times d}$  denote the *expected  $n$ -sample Gramian matrix* given by

$$\forall i, i' \in [d], [\Sigma_n]_{i, i'} = \frac{1}{n} \sum_{j=1}^n \mathbb{E}_{X^{(j)} \sim P_X^{(j)}} [\langle [\Phi_j]_i, [\Phi_j]_{i'} \rangle] \quad \therefore \Sigma_n = \mathbb{E}_{X^{1:n}} \left[ \frac{1}{n} \sum_{j=1}^n \int_S \Phi_j \Phi_j^\top d\mathbf{m} \right].$$

As mentioned previously, the objective of contextual CDF regression is to recover  $\theta_*$  given the samples  $(x, y)^{1:n}$ . Each  $Y^{(j)}$  defines a one-sample empirical CDF  $\mathbf{I}_{Y^{(j)}}(t) = \mathbb{1}\{t \geq Y^{(j)}\}$  which approximates the true CDF  $F(x^{(j)}, \cdot) = \mathbb{E}[\mathbf{I}_{Y^{(j)}}(\cdot) \mid X^{(j)} = x^{(j)}]$  in expectation conditioned on  $X^{(j)} = x^{(j)}$ . We investigate three estimators for  $\theta_*$  based on the paradigm of empirical risk minimization. Firstly, the *lasso estimator*  $\hat{\theta}_\lambda$  imposes an  $\ell^1$ -penalty on the parameter vector, weighted by  $\lambda > 0$ :

$$\hat{\theta}_\lambda = \arg \min_{\theta \in \mathbb{R}^d} \left\{ \frac{1}{n} \sum_{j=1}^n \|\mathbf{I}_{y^{(j)}} - \theta^\top \Phi_j\|_{\mathcal{L}^2(S, \mathbf{m})}^2 + \lambda \|\theta\|_1 \right\}. \quad (1)$$

Secondly, the *elastic net estimator*  $\hat{\theta}_{\lambda_1, \lambda_2}$  imposes both  $\ell^1$ - and  $\ell^2$ -penalties on the parameter vector, weighted by  $\lambda_1 > 0$  and  $\lambda_2 > 0$  respectively:

$$\hat{\theta}_{\lambda_1, \lambda_2} = \arg \min_{\theta \in \mathbb{R}^d} \left\{ \frac{1}{n} \sum_{j=1}^n \|\mathbf{I}_{y^{(j)}} - \theta^\top \Phi_j\|_{\mathcal{L}^2(S, \mathbf{m})}^2 + \lambda_1 \|\theta\|_1 + \lambda_2 \|\theta\|^2 \right\}. \quad (2)$$

Thirdly, the *Dantzig selector*  $\bar{\theta}_\lambda$  is given by the following variation on the constrained optimization formulation of lasso estimation:

$$\bar{\theta}_\lambda = \underset{\theta \in \mathbb{R}^d}{\text{minimize}} \quad \|\theta\|_1 \quad (3)$$

$$\text{subject to} \quad \left\| \frac{1}{n} \sum_{j=1}^n \langle \mathbf{I}_{y^{(j)}} - \theta^\top \Phi_j, \Phi_j \rangle \right\|_\infty \leq \lambda. \quad (4)$$

For all three estimators, we analyze the *estimation error* given by the Euclidean distance  $\|\hat{\theta} - \theta_*\|_A$  between the estimated and true parameter vectors, possibly weighted by some matrix  $A$ . We remark that Equations (1) to (3) define *improper estimators* which do not necessarily lie in  $\Delta^{d-1}$  and are thus not guaranteed to define valid CDFs. However, since  $\theta_* \in \Delta^{d-1}$  and  $\Delta^{d-1}$  is closed and convex, any upper bound on the estimation error of an improper estimator also holds for its projection onto  $\Delta^{d-1}$ , by Beck (2014, Equation 9.10).

Lastly, we remark that our results hold for any  $s \leq d$ , although we are principally interested in the advantages of our proposed methods over the ridge estimator baseline (Zhang et al., 2024, Section 3.1) for small values of  $s$ . We emphasize that our results have no dependence on the cardinality of the context space  $\mathcal{X}$ .

## 5 Main Results

### 5.1 Lasso

Before presenting our main contributions, we introduce the following condition on symmetric matrices which is used in the statements of our main results:

**Definition 1** (Restricted eigenvalue condition). *Fix  $d \in \mathbb{N}$ ,  $\kappa \geq 0$ ,  $\gamma \geq 0$ , and  $S_* \subseteq [d]$ . Let*

$$C_\gamma(S_*) = \left\{ v \in \mathbb{R}^d : \left\| [v]_{S_*^c} \right\|_1 \leq \gamma \left\| [v]_{S_*} \right\|_1 \right\}.$$

*A symmetric matrix  $A \in \mathbb{R}^{d \times d}$  satisfies the  $(\kappa, \gamma)$ -restricted eigenvalue condition over  $S_*$  iff*

$$\forall v \in C_\gamma(S_*), \quad v^\top A v \geq \kappa \|v\|^2.$$

Intuitively, the restricted eigenvalue condition specializes the notions of positive definiteness and strong convexity to the subset of directions whose support is close to  $S_*$ . The significance of this condition lies in its application to the Gramian matrices  $U_n$  and  $\Sigma_n$ , since lower-bounding the Gramian eigenvalues is sufficient to ensure the well-conditioning of the basis functions comprising the true CDF. Our formulation of the restricted eigenvalue condition extends the standard definition from Wainwright (2019, Definition 7.12) to the general inner product spaces used in the definitions of  $U_n$  and  $\Sigma_n$ . We provide examples of non-trivial CDF bases which satisfy the restricted eigenvalue condition in Appendix F.

Under this assumption, we state our first main result, a high-probability upper bound on the error of the lasso estimator in the fixed design setting:

**Theorem 1** (Lasso fixed design upper bound). *Fix  $\delta \in (0, 1)$ . Assume the samples  $(x, y)^{1:n}$  are generated according to the fixed design setting. Assume  $U_n$  satisfies the  $(\kappa, 3)$ -restricted eigenvalue condition over  $S_*$ . Let  $\hat{\theta}_\lambda$  be the lasso estimator (1) with regularization hyperparameter  $\lambda = 4\sqrt{(2/n)\log(2d/\delta)}$ . Then with probability at least  $1 - \delta$ , the estimation error satisfies the bounds*

$$\left\| \hat{\theta}_\lambda - \theta_* \right\|_{U_n} \leq 6\sqrt{\frac{2s}{\kappa n} \log\left(\frac{2d}{\delta}\right)} \quad \text{and} \quad \left\| \hat{\theta}_\lambda - \theta_* \right\| \leq \frac{6}{\kappa} \sqrt{\frac{2s}{n} \log\left(\frac{2d}{\delta}\right)}.$$

We provide the technical details of our proof in Section 6. The crux of our argument lies in the definition of  $\hat{\theta}_\lambda$  as the minimizer of the empirical risk objective in (1). Substituting  $\hat{\theta}_\lambda$  and  $\theta_*$  into this objective results in an inequality (10) in terms of  $\hat{\theta}_\lambda$  and  $\theta_*$ . Rearranging yields an upper-bound (12) on the estimation error  $\|\Delta\|_{U_n}$  in terms of two quantities:

- The scalar projections of the sampling errors  $\mathbf{I}_{y^{(j)}} - \theta_*^\top \Phi_j$  onto the  $d$  basis functions  $\Phi_j$ , which we upper-bound with high probability (7) using various concentration inequalities.
- The difference between the  $\ell^1$ -regularizers  $\|\theta_*\|_1 - \|\hat{\theta}_\lambda\|_1$ , which we upper-bound by utilizing the sparsity of  $\theta_*$  and the definition of  $\ell^1$ -norm, among other techniques.

Combining these results gives rise to a bound on the relative magnitudes of the components of  $\Delta$  corresponding to the true support  $S_*$  and its complement (15). Then, upper-bounding the component  $\|[\Delta]_{S_*}\|_1$  is sufficient (16) to characterize the overall estimation error  $\|\Delta\|_{U_n}$ . Applying various norm equivalences and invoking the restricted eigenvalue condition on  $U_n$  completes the proof.

Theorem 1 establishes the asymptotic complexity of lasso CDF regression as  $O(\sqrt{s \log(d)/n}) = \tilde{O}(\sqrt{s/n})$ , analogously to the  $\tilde{O}(\sqrt{d/n})$  result for ridge regression in Zhang et al. (2024). The dependence on  $\sqrt{s}$  arises from our sparsity analysis on  $\|\theta_*\|_1 - \|\hat{\theta}_\lambda\|_1$  and the upper-bound on  $\|[\Delta]_{S_*}\|_1$  as discussed above. The factors of  $1/\sqrt{n}$  and  $\sqrt{\log(d)}$  arise from the probabilistic arguments in our analysis of the sampling errors. Ultimately, the lasso estimator's characteristic  $\ell^1$ -penalty term plays an instrumental role in deriving a result with sub-polynomial dependence on  $d$ .

Our second main result is a high-probability upper bound on the error of the lasso estimator in the random design setting, under similar assumptions on the well-conditioning of the CDF basis:

**Theorem 2** (Lasso random design upper bound). *Fix  $\delta_1 \in (0, 1)$  and  $\delta_2 \in (0, 1 - \delta_1)$ . Assume the samples  $(x, y)^{1:n}$  are generated according to the random design setting. Assume that for any  $n \in \mathbb{N}$ , the matrix  $\Sigma_n$  satisfies the  $(\kappa, 3)$ -restricted eigenvalue condition over  $S_*$ . Assume the sample size is at least  $n \geq (32d^2/\kappa^2) \log(d/\delta_1)$ . Let  $\hat{\theta}_\lambda$  be the lasso estimator (1) with regularization hyperparameter  $\lambda = 4\sqrt{(2/n) \log(2d/\delta_2)}$ . Then with probability at least  $1 - \delta_1 - \delta_2$ , the estimation error satisfies the bounds*

$$\left\| \hat{\theta}_\lambda - \theta_* \right\|_{U_n} \leq 12 \sqrt{\frac{s}{\kappa n} \log\left(\frac{2d}{\delta_2}\right)} \quad \text{and} \quad \left\| \hat{\theta}_\lambda - \theta_* \right\| \leq \frac{12}{\kappa} \sqrt{\frac{2s}{n} \log\left(\frac{2d}{\delta_2}\right)}.$$

We defer the technical details of our proof to Appendix B. In a nutshell, since  $U_n$  is a sum of  $n$  independent matrices and is an unbiased estimator of  $\Sigma_n$ , we can employ a matrix analogue of Hoeffding's inequality (17) to justify approximating  $\Sigma_n$  with  $U_n$  given sufficiently many samples  $n$ . Subsequently, we analyze the estimation error weighted by  $U_n$  to complete the proof.

## 5.2 Elastic Net

Our third main result is a high-probability upper bound on the error of the elastic net estimator in the fixed design setting:

**Theorem 3** (Elastic net upper bound). *Fix  $\delta \in (0, 1)$ . Assume the samples  $(x, y)^{1:n}$  are generated according to the fixed design setting. Assume  $U_n$  satisfies the  $(\kappa, 3 + 4\lambda_2/\lambda_1)$ -restricted eigenvalue condition over  $S_*$ . Let  $\hat{\theta}_{\lambda_1, \lambda_2}$  be the elastic net estimator (2) with  $\ell^1$ -regularization hyperparameter  $\lambda_1 = 4\sqrt{(2/n) \log(2d/\delta)}$ . Then with probability at least  $1 - \delta$ , the estimation error satisfies the bounds*

$$\begin{aligned} \left\| \hat{\theta}_{\lambda_1, \lambda_2} - \theta_* \right\|_{U_n + \lambda_2 I_d} &\leq \left( 6 \sqrt{\frac{2}{n} \log\left(\frac{2d}{\delta}\right)} + 2\lambda_2 \right) \sqrt{\frac{s}{\kappa + \lambda_2}}, \\ \left\| \hat{\theta}_{\lambda_1, \lambda_2} - \theta_* \right\| &\leq \left( 6 \sqrt{\frac{2}{n} \log\left(\frac{2d}{\delta}\right)} + 2\lambda_2 \right) \frac{\sqrt{s}}{\kappa + \lambda_2}. \end{aligned} \tag{5}$$

Furthermore, if the  $\ell^2$ -regularization hyperparameter is  $\lambda_2 = 3\sqrt{(2/n)\log(2d/\delta)}$ , the bound (5) implies

$$\left\| \hat{\theta}_{\lambda_1, \lambda_2} - \theta_* \right\|_{U_n + \lambda_2 I_d} \leq 4 \sqrt{\frac{18s^2}{n} \log\left(\frac{2d}{\delta}\right)}. \quad (6)$$

We emphasize that Theorem 3 produces non-trivial estimation bounds even when no assumptions are placed on  $U_n$  (i.e.,  $\kappa = 0$ ). At a high level, the elastic net estimator with its characteristic  $\ell^2$ -penalty term may be perceived as “regularizing” the empirical Gramian matrix to produce  $U_n + \lambda_2 I_d$ , whose smallest eigenvalue is strictly positive. Additional conditions on  $U_n$  further tighten the estimation bounds by bolstering the minimum eigenvalue of  $U_n + \lambda_2 I_d$ . Hence, the elastic net estimator simultaneously selects features and regularizes the problem in an integrated fashion by ensuring strong convexity of the objective function (2).

We defer the technical details of our proof to Appendix C. We carry out an analysis of the sampling errors and upper-bound the difference between the  $\ell^1$ -regularizers  $\|\theta_*\|_1 - \|\hat{\theta}_{\lambda_1, \lambda_2}\|_1$ , with additional accommodations for the elastic net estimator’s  $\ell^2$ -penalty term. Notably, substituting  $\hat{\theta}_{\lambda_1, \lambda_2}$  and  $\theta_*$  into the empirical risk objective (2) and rearranging yields an upper-bound on  $\|\Delta\|_{U_n + \lambda_2 I_d}$  instead of  $\|\Delta\|_{U_n}$ . The resulting expression (22) contains an additional quantity  $\lambda_2 \Delta^\top \theta_*$ , which we upper-bound in terms of the error component  $\|[\Delta]_{S_*}\|_1$  in (23). Combining this contribution with the terms resulting from the  $\ell^1$ -regularizers gives rise to the parenthesized sum in (5), and choosing the optimal value of  $\lambda_2$  to balance the summands achieves the bound in (6).

### 5.3 Dantzig Selector

Before presenting our fourth main result, we introduce two more conditions on symmetric matrices which underlie standard assumptions in the statistics and compressed sensing literature (Bandeira et al., 2013):

**Definition 2** (Restricted isometry property). *Fix  $\epsilon \geq 0$  and  $p \in \mathbb{N}$ . A symmetric matrix  $A \in \mathbb{R}^{d \times d}$  satisfies the  $(\epsilon, p)$ -restricted isometry property iff*

$$\forall v \in \mathbb{R}^d, \|v\|_0 \leq p \implies (1 - \epsilon) \|v\|^2 \leq v^\top A v \leq (1 + \epsilon) \|v\|^2.$$

The restricted isometry property states that the curvature of the quadratic form represented by  $A$  is bounded around 1 along directions residing in sparse axis-aligned subspaces, or alternatively, that the linear operator represented by  $A$  is approximately scale-preserving when applied to sparse inputs. A weaker variant of this property bears resemblance to the Cauchy-Schwarz inequality:

**Definition 3** (Restricted orthogonality property). *Fix  $\zeta \geq 0$ ,  $p \in \mathbb{N}$ , and  $q \in \mathbb{N}$ . A symmetric matrix  $A \in \mathbb{R}^{d \times d}$  satisfies the  $(\zeta, p, q)$ -restricted orthogonality property iff*

$$\forall u, v \in \mathbb{R}^d, \|u\|_0 \leq p \wedge \|v\|_0 \leq q \wedge \text{supp}(u) \cap \text{supp}(v) = \emptyset \implies |u^\top A v| \leq \zeta \|u\| \|v\|.$$

Under the restricted orthogonality property,  $A$  maps sparse vectors with disjoint support to dissimilar outputs. Our fourth main result is a high-probability upper bound on the error of the Dantzig selector in the fixed design setting, under assumptions on  $U_n$  similar to the prior literature (Candes & Tao, 2007):

**Theorem 4** (Dantzig selector upper bound). *Fix  $\delta \in (0, 1)$ ,  $\epsilon \in [0, 1)$ , and  $\zeta \in [0, 1 - \epsilon)$ . Assume the samples  $(x, y)^{1:n}$  are generated according to the fixed design setting. Assume  $U_n$  satisfies the  $(\epsilon, 2s)$ -restricted isometry property and the  $(\zeta, s, 2s)$ -restricted orthogonality property. Let  $\bar{\theta}_\lambda$  be the Dantzig selector (3) with regularization hyperparameter  $\lambda = \sqrt{(2/n)\log(2d/\delta)}$ . Then with probability at least  $1 - \delta$ , the estimation error satisfies the bound*

$$\|\bar{\theta}_\lambda - \theta_*\| \leq \frac{4}{1 - \epsilon - \zeta} \sqrt{\frac{2s}{n} \log\left(\frac{2d}{\delta}\right)}.$$

We include the Dantzig selector in our investigation of sparse contextual CDF regression as an example of variable selection formulated as a linear programming problem (4), as opposed to the usual quadratic



programming perspective of lasso estimation. We remark that the  $(\epsilon, 3p)$ -restricted isometry property implies the  $(\epsilon, p, 2p)$ -restricted orthogonality property, by [Candes & Tao \(2005, Lemma 1.1\)](#). Hence, [Theorem 4](#) also holds when the sole assumption placed on  $U_n$  is the  $(\epsilon, 3s)$ -restricted isometry property.

We defer the technical details of our proof to [Appendix D](#) and provide a high-level summary below. At the outset, we upper-bound the estimation error  $\|\Delta\|$  in terms of the components  $\|[\Delta]_{S_\dagger}\|$  and  $\|[\Delta]_{S_\dagger^c}\|_1$ , where  $S_\dagger$  is a superset of  $S_*$  ([30](#)). Utilizing the sparsity of  $\theta_*$ , we bound the latter component in terms of the former ([31](#)). Subsequently, we upper-bound  $\|[\Delta]_{S_\dagger}\|$  in terms of  $\|[U_n]_{\langle S_\dagger \rangle} \Delta\|$  ([33](#)), which in turn depends on two quantities ([34](#)):

- The scalar projections of the sampling errors  $I_{y^{(j)}} - \theta_*^\top \Phi_j$  onto the  $d$  basis functions  $\Phi_j$ , which we upper-bound using various probabilistic arguments and concentration inequalities.
- The scalar projections of  $I_{y^{(j)}} - \bar{\theta}_\lambda^\top \Phi_j$  onto  $\Phi_j$ , which we upper-bound using the constraint in ([4](#)).

Combining these bounds and applying various norm equivalences completes the proof. We note that the restricted isometry and restricted orthogonality assumptions on  $U_n$  are preconditions to upper-bound  $\|\Delta\|$  and  $\|[\Delta]_{S_\dagger}\|$  in the proof.

## 5.4 Lower Bound

Our last main result is a lower bound on the *minimax*  $\ell^2$ -risk of sparse contextual CDF regression. To begin this discussion, we introduce some relevant notation. For any  $d \in \mathbb{N}$ , let  $\mathcal{B}_d$  be the universe of all  $d$ -dimensional bases of contextual CDFs, i.e.,

$$\mathcal{B}_d = \left\{ \Phi : \mathcal{X} \times \mathbb{R} \rightarrow [0, 1]^d : \forall i \in [d], \forall x \in \mathcal{X}, [\Phi(x, \cdot)]_i \text{ is a CDF} \right\}.$$

For any  $x \in \mathcal{X}$ ,  $\theta \in \Delta^{d-1}$ , and  $\Phi \in \mathcal{B}_d$ , let  $P_{Y|x, \theta}^\Phi$  denote the probability distribution corresponding to the CDF  $\theta^\top \Phi(x, \cdot)$ . Given context variables  $x^{1:n} \in \mathcal{X}^n$ , let  $\mathcal{P}_{x^{1:n}}^d$  denote the family of product distributions of  $Y^{1:n}$  which are convex combinations of  $d$  contextual CDFs, i.e.,

$$\mathcal{P}_{x^{1:n}}^d = \left\{ \bigotimes_{j=1}^n P_{Y|x^{(j)}, \theta}^\Phi : \theta \in \Delta^{d-1} \wedge \Phi \in \mathcal{B}_d \right\}.$$

For any  $s \leq d$ , let  $\mathcal{P}_{x^{1:n}}^{d,s} \subset \mathcal{P}_{x^{1:n}}^d$  denote the family of product distributions of  $Y^{1:n}$  which are convex combinations of  $s$  contextual CDFs chosen from a  $d$ -dimensional basis, i.e.,

$$\mathcal{P}_{x^{1:n}}^{d,s} = \left\{ \bigotimes_{j=1}^n P_{Y|x^{(j)}, \theta}^\Phi : \theta \in \Delta^{d-1} \wedge \|\theta\|_0 = s \wedge \Phi \in \mathcal{B}_d \right\}.$$

Given a distribution  $P \in \mathcal{P}_{x^{1:n}}^d$ , let  $\theta(P) \in \Delta^{d-1}$  denote its parameter and let  $S_*(P) = \text{supp}(\theta(P)) \subseteq [d]$ . Consequently,  $|S_*(P)| = s$  for any  $P \in \mathcal{P}_{x^{1:n}}^{d,s}$ . For any  $n \in \mathbb{N}$ , let  $\hat{\Theta}_{d,n}$  be the universe of all (possibly randomized) estimators  $\hat{\theta} : \mathbb{R}^n \rightarrow \mathbb{R}^d$ .

Our main result establishes a lower bound on the estimation error of any estimator for sparse contextual CDF regression, and is obtained as a consequence of the lower bound for general contextual CDF regression in [Zhang et al. \(2024, Theorem 8\)](#):

**Proposition 1** (Lower bound). *Fix  $d \in \mathbb{N}$ ,  $s \leq d$ , and any sufficiently large  $n \geq s/2$ . Fix  $x^{1:n} \in \mathcal{X}^n$ . Then, the minimax  $\ell^2$ -risk of sparse contextual CDF regression satisfies the bound*

$$\mathfrak{R}\left(\theta\left(\mathcal{P}_{x^{1:n}}^{d,s}\right)\right) = \inf_{\hat{\theta} \in \hat{\Theta}_{d,n}} \sup_{P \in \mathcal{P}_{x^{1:n}}^{d,s}} \mathbb{E}_{Y^{1:n} \sim P} \left[ \left\| \hat{\theta}(Y^{1:n}) - \theta(P) \right\| \right] = \Omega\left(\sqrt{\frac{s}{n}}\right).$$

We defer the technical details of our proof to Appendix E. In a nutshell, it suffices to consider only the estimation error component corresponding to the indices in  $S_*(P)$ , thereby reducing the problem to general  $s$ -dimensional CDF regression. Furthermore, we obtain minimax upper bounds from the high-probability upper bounds in Theorems 1 and 4 by the arguments from Zhang et al. (2024, p. 12). Hence, Proposition 1 establishes the minimax optimality of the lasso and Dantzig selector estimators up to logarithmic factors.

## 6 Proof of Lasso Fixed Design Upper Bound

In this section, we prove Theorem 1. We begin by deriving a concentration bound on the inner products between the sampling errors  $I_{y^{(j)}} - \theta_*^\top \Phi_j$  and the basis functions  $\Phi_j$ :

**Lemma 1** (Concentration bound). *Fix  $\delta \in (0, 1)$ . With probability at least  $1 - \delta$ , it holds that*

$$\left\| \frac{1}{n} \sum_{j=1}^n \langle I_{y^{(j)}} - \theta_*^\top \Phi_j, \Phi_j \rangle \right\|_\infty \leq \sqrt{\frac{2}{n} \log\left(\frac{2d}{\delta}\right)}. \quad (7)$$

*Proof of Lemma 1.* For each  $j \in [n]$  and  $i \in [d]$ , define a random variable

$$Z_{i,j} = \langle I_{Y^{(j)}} - \theta_*^\top \Phi_j, [\Phi_j]_i \rangle.$$

Let  $\mathcal{F}_j = \sigma(Y^{1:j-1})$  denote the  $\sigma$ -algebra generated by the random variables  $Y^{1:j-1}$ .<sup>1</sup> The mean of  $Z_{i,j}$  is

$$\begin{aligned} \mathbb{E}[Z_{i,j}] &\stackrel{(a)}{=} \mathbb{E}[\mathbb{E}[Z_{i,j} \mid \mathcal{F}_j]] \\ &\stackrel{(b)}{=} \mathbb{E}[\mathbb{E}[\langle I_{Y^{(j)}} - \theta_*^\top \Phi_j, [\Phi_j]_i \rangle \mid \mathcal{F}_j]] \\ &\stackrel{(c)}{=} \mathbb{E}[\langle \mathbb{E}[I_{Y^{(j)}} - \theta_*^\top \Phi_j \mid \mathcal{F}_j], [\Phi_j]_i \rangle] \\ &\stackrel{(d)}{=} \mathbb{E}[\langle 0, [\Phi_j]_i \rangle] = 0, \end{aligned}$$

where (a) holds by the tower rule of expectation, (b) holds by definition of  $Z_{i,j}$ , (c) holds by Fubini's theorem, and (d) holds because  $\mathbb{E}[I_{Y^{(j)}} \mid \mathcal{F}_j] = \theta_*^\top \Phi_j$ . Furthermore, the support of  $Z_{i,j}$  is  $[0, 1]$  because

$$\begin{aligned} |Z_{i,j}| &\stackrel{(a)}{=} \left| \int_S (I_{y^{(j)}}(t) - \theta_*^\top \Phi_j(t)) \phi_i(x^{(j)}, t) d\mathbf{m} \right| \\ &\stackrel{(b)}{\leq} \int_S |I_{y^{(j)}}(t) - \theta_*^\top \Phi_j(t)| |\phi_i(x^{(j)}, t)| d\mathbf{m} \\ &\stackrel{(c)}{\leq} \mathbf{m}(S) \stackrel{(d)}{=} 1, \end{aligned}$$

where (a) holds by definition of inner product, (b) holds by the triangle inequality, (c) holds because CDFs are bounded between 0 and 1, and (d) holds because  $\mathbf{m}$  is a probability measure. Thus, for any  $\tau > 0$ ,

$$\mathbb{P}\left(\left\| \frac{1}{n} \sum_{j=1}^n \langle I_{Y^{(j)}} - \theta_*^\top \Phi_j, \Phi_j \rangle \right\|_\infty \geq \tau\right) \stackrel{(a)}{\leq} \sum_{i=1}^d \mathbb{P}\left(\left| \frac{1}{n} \sum_{j=1}^n Z_{i,j} \right| \geq \tau\right) \stackrel{(b)}{\leq} 2d \exp\left(-\frac{n\tau^2}{2}\right),$$

where (a) holds by definition of  $Z_{i,j}$  and the union bound, and (b) holds by Hoeffding's inequality. Choosing  $\tau = \sqrt{2/n \log(2d/\delta)}$  and rearranging, we get  $\delta = 2d \exp(-n\tau^2/2)$ , and thus

$$\mathbb{P}\left(\left\| \frac{1}{n} \sum_{j=1}^n \langle I_{Y^{(j)}} - \theta_*^\top \Phi_j, \Phi_j \rangle \right\|_\infty \leq \sqrt{\frac{2}{n} \log\left(\frac{2d}{\delta}\right)}\right) \geq 1 - \delta$$

as desired. ■

<sup>1</sup>This lemma also applies in the random design setting by taking  $\mathcal{F}_j = \sigma(X^{1:j}, Y^{1:j-1})$ .



Now, we are ready to prove Theorem 1.

*Proof of Theorem 1.* Throughout this proof, we restrict to the subset of the probability space where

$$\left\| \frac{1}{n} \sum_{j=1}^n \langle \mathbf{I}_{y^{(j)}} - \theta_*^\top \Phi_j, \Phi_j \rangle \right\|_\infty \leq \sqrt{\frac{2}{n} \log \left( \frac{2d}{\delta} \right)} = \frac{\lambda}{4}, \quad (8)$$

which holds with probability at least  $1 - \delta$  by Lemma 1. For notational simplicity, let  $\Delta = \hat{\theta}_\lambda - \theta_*$ . We have

$$\begin{aligned} \|\Delta\|_{U_n}^2 &\stackrel{(a)}{=} \Delta^\top U_n \Delta \\ &\stackrel{(b)}{=} \Delta^\top \left( \frac{1}{n} \sum_{j=1}^n \int_S \Phi_j \Phi_j^\top d\mathbf{m} \right) \Delta \\ &\stackrel{(c)}{=} \frac{1}{n} \sum_{j=1}^n \int_S \Delta^\top \Phi_j \Phi_j^\top \Delta d\mathbf{m} \\ &\stackrel{(d)}{=} \frac{1}{n} \sum_{j=1}^n \|\Delta^\top \Phi_j\|_{\mathcal{L}^2(S, \mathbf{m})}^2 \\ &\stackrel{(e)}{=} \frac{1}{n} \sum_{j=1}^n \left( \|\hat{\theta}_\lambda^\top \Phi_j\|_{\mathcal{L}^2(S, \mathbf{m})}^2 + \|\theta_*^\top \Phi_j\|_{\mathcal{L}^2(S, \mathbf{m})}^2 \right) - \frac{2}{n} \sum_{j=1}^n \langle \hat{\theta}_\lambda^\top \Phi_j, \theta_*^\top \Phi_j \rangle \\ &\stackrel{(f)}{=} \frac{1}{n} \sum_{j=1}^n \left( \|\hat{\theta}_\lambda^\top \Phi_j\|_{\mathcal{L}^2(S, \mathbf{m})}^2 - \|\theta_*^\top \Phi_j\|_{\mathcal{L}^2(S, \mathbf{m})}^2 \right) - \frac{2}{n} \sum_{j=1}^n \langle \Delta^\top \Phi_j, \theta_*^\top \Phi_j \rangle, \end{aligned} \quad (9)$$

where (a) holds by definition of weighted  $\ell^2$ -norm induced by  $U_n$ , (b) holds by definition of  $U_n$ , (c) holds by the linearity of integration, (d) holds by definition of inner product between functions, (e) holds by definition of  $\Delta$ , and (f) holds by definition of  $\Delta$  and the linearity of inner product. Next, since  $\hat{\theta}_\lambda$  minimizes the objective in (1), it follows that

$$\frac{1}{n} \sum_{j=1}^n \left\| \mathbf{I}_{y^{(j)}} - \hat{\theta}_\lambda^\top \Phi_j \right\|_{\mathcal{L}^2(S, \mathbf{m})}^2 + \lambda \|\hat{\theta}_\lambda\|_1 \leq \frac{1}{n} \sum_{j=1}^n \left\| \mathbf{I}_{y^{(j)}} - \theta_*^\top \Phi_j \right\|_{\mathcal{L}^2(S, \mathbf{m})}^2 + \lambda \|\theta_*\|_1. \quad (10)$$

Expanding the squared norms and rearranging, it follows that (cf. Wainwright (2019))

$$\frac{1}{n} \sum_{j=1}^n \left( \|\hat{\theta}_\lambda^\top \Phi_j\|_{\mathcal{L}^2(S, \mathbf{m})}^2 - \|\theta_*^\top \Phi_j\|_{\mathcal{L}^2(S, \mathbf{m})}^2 \right) \leq \frac{2}{n} \sum_{j=1}^n \langle \mathbf{I}_{y^{(j)}}, \Delta^\top \Phi_j \rangle + \lambda \left( \|\theta_*\|_1 - \|\hat{\theta}_\lambda\|_1 \right). \quad (11)$$

Combining Equations (9) and (11), we obtain

$$\|\Delta\|_{U_n}^2 \leq \underbrace{\frac{2}{n} \sum_{j=1}^n \langle \mathbf{I}_{y^{(j)}}, \Delta^\top \Phi_j \rangle - \frac{2}{n} \sum_{j=1}^n \langle \Delta^\top \Phi_j, \theta_*^\top \Phi_j \rangle}_{\textcircled{1}} + \underbrace{\lambda \left( \|\theta_*\|_1 - \|\hat{\theta}_\lambda\|_1 \right)}_{\textcircled{2}}. \quad (12)$$

Next, we upper-bound  $\textcircled{1}$ . We have

$$\textcircled{1} \stackrel{(a)}{=} \Delta^\top \left( \frac{2}{n} \sum_{j=1}^n \langle \mathbf{I}_{y^{(j)}} - \theta_*^\top \Phi_j, \Phi_j \rangle \right) \stackrel{(b)}{\leq} \|\Delta\|_1 \left\| \frac{2}{n} \sum_{j=1}^n \langle \mathbf{I}_{y^{(j)}} - \theta_*^\top \Phi_j, \Phi_j \rangle \right\|_\infty \stackrel{(c)}{\leq} \frac{\lambda}{2} \|\Delta\|_1, \quad (13)$$

where (a) holds by the linearity of inner product, (b) follows from Hölder's inequality, and (c) holds by substituting in (8). Next, we upper-bound  $\textcircled{2}$ . We have

$$\textcircled{2} \stackrel{(a)}{=} \sum_{i=1}^d |[\theta_*]_i| - \sum_{i=1}^d |[\theta_* + \Delta]_i|$$

$$\begin{aligned}
&\stackrel{(b)}{=} \sum_{i \in S_*} |[\theta_*]_i| - \sum_{i \in S_*} |[\theta_* + \Delta]_i| - \sum_{i \in S_*^c} |[\Delta]_i| \\
&\stackrel{(c)}{=} \|[\theta_*]_{S_*}\|_1 - \|[\theta_*]_{S_*} + [\Delta]_{S_*}\|_1 - \|[\Delta]_{S_*^c}\|_1 \\
&\stackrel{(d)}{\leq} \|[\Delta]_{S_*}\|_1 - \|[\Delta]_{S_*^c}\|_1,
\end{aligned} \tag{14}$$

where (a) holds by definition of  $\ell^1$ -norm and  $\Delta$ , (b) holds because  $\text{supp}(\theta_*) = S_*$ , (c) holds by definition of  $\ell^1$ -norm, and (d) holds by the triangle inequality. Combining Equations (12) to (14), we obtain

$$\|\Delta\|_{U_n}^2 \leq \lambda \left( \frac{3}{2} \|[\Delta]_{S_*}\|_1 - \frac{1}{2} \|[\Delta]_{S_*^c}\|_1 \right) \tag{15}$$

$$\leq \frac{3\lambda}{2} \|[\Delta]_{S_*}\|_1 \tag{16}$$

$$\stackrel{(a)}{\leq} \frac{3\lambda}{2} \sqrt{s} \|[\Delta]_{S_*}\| \leq \frac{3\lambda}{2} \sqrt{s} \|\Delta\|$$

$$\stackrel{(b)}{\leq} \frac{3\lambda}{2} \sqrt{\frac{s}{\kappa}} \|\Delta\|_{U_n},$$

where (a) holds by the equivalence between  $\ell^1$ - and  $\ell^2$ -norms, and (b) follows from the restricted eigenvalue condition on  $U_n$  since rearranging (15) yields  $\|[\Delta]_{S_*^c}\|_1 \leq 3\|[\Delta]_{S_*}\|_1$ . Thus, the estimation error satisfies the bounds

$$\|\Delta\|_{U_n} \leq \frac{3\lambda}{2} \sqrt{\frac{s}{\kappa}} \stackrel{(a)}{=} 6\sqrt{\frac{2s}{\kappa n} \log\left(\frac{2d}{\delta}\right)} \quad \text{and} \quad \|\Delta\| \leq \frac{1}{\sqrt{\kappa}} \|\Delta\|_{U_n} = \frac{6}{\kappa} \sqrt{\frac{2s}{n} \log\left(\frac{2d}{\delta}\right)},$$

where (a) holds by substituting in  $\lambda$  and (b) follows from the restricted eigenvalue condition on  $U_n$ .  $\blacksquare$

## 7 Numerical Simulations

In this section, we numerically simulate the data generation process described in Section 4 and empirically evaluate the accuracy of our proposed lasso estimator (1) on synthetic data. For mathematical conciseness, we consider a basis of contextual Bernoulli CDFs, which yields a closed-form expression for the induced norm  $\|\cdot\|_{\mathcal{L}^2(S, \mathbf{m})}$  in the training objective. Formally, let  $\mathcal{X} = [0, 1]^d$  be the space of all  $d$ -tuples  $x = (x_1, \dots, x_d)$  of Bernoulli parameters, and define the basis functions

$$\forall i \in [d], \phi_i(x, t) = \begin{cases} 1 - x_i, & \text{if } 0 \leq t < 1, \\ 1, & \text{if } t = 1, \end{cases}$$

where  $\phi_i(x, \cdot)$  is the CDF of a  $\text{Bernoulli}(x_i)$  random variable. It follows that  $F(x, \cdot) = \theta_*^\top \Phi(x, \cdot)$  is the CDF of a  $\text{Bernoulli}(\theta_*^\top x)$  random variable, because

$$F(x, t) = \sum_{i=1}^d [\theta_*]_i \cdot \begin{cases} 1 - x_i, & \text{if } 0 \leq t < 1 \\ 1, & \text{if } t = 1 \end{cases} \stackrel{(a)}{=} \begin{cases} 1 - \theta_*^\top x, & \text{if } 0 \leq t < 1, \\ 1, & \text{if } t = 1, \end{cases}$$

where (a) holds because  $\theta_*$  is a PMF. Let  $\mathbf{m}$  be the uniform measure over the support  $S = [0, 1]$ . Then, the log-likelihood term in (1) is equivalent to the canonical least squares formulation because

$$\begin{aligned}
\sum_{j=1}^n \|\mathbf{I}_{y^{(j)}} - \theta^\top \Phi_j\|_{\mathcal{L}^2(S, \mathbf{m})}^2 &\stackrel{(a)}{=} \sum_{j=1}^n \int_0^1 (\mathbf{I}_{y^{(j)}}(t) - \theta^\top \Phi_j(t))^2 dt \\
&\stackrel{(b)}{=} \sum_{j=1}^n \left( \mathbb{1}\{y^{(j)} = 0\} - (1 - \theta^\top x^{(j)}) \right)^2
\end{aligned}$$

$$\begin{aligned}
&= \sum_{j=1}^n \left( \theta^\top x^{(j)} - y^{(j)} \right)^2 \\
&\stackrel{(c)}{=} \|A\theta - b\|^2,
\end{aligned}$$

where (a) follows from the support  $S = [0, 1]$  and the choice of measure, (b) holds by substituting in the definitions of  $\mathbf{I}_{y^{(j)}}$  and  $\Phi_j$ , and the matrix  $A \in \mathbb{R}^{n \times d}$  and vector  $b \in \mathbb{R}^n$  in (c) are defined as  $[A]_{\langle j \rangle} = x^{(j)\top}$  and  $[b]_j = y^{(j)}$ . We consider the fixed design setting with context variables  $x^{(j)} = (x_1^{(j)}, \dots, x_d^{(j)})$  given by

$$\begin{aligned}
x_i^{(j)} &= \begin{cases} 1 - 2x_{\text{val}}^{(j)}, & \text{if } i \equiv j \pmod{d} \\ 1 - x_{\text{val}}^{(j)}, & \text{if } i \not\equiv j \pmod{d} \end{cases} \quad \text{with} \quad x_{\text{val}}^{(j)} = \begin{cases} \frac{1}{2}, & \text{if } j \leq d \\ \frac{\mu_{\min}(M_{j-1})}{\alpha_j}, & \text{if } j > d \end{cases} \\
\text{and } M_j &= \left(1 - x^{(j)}\right) \left(1 - x^{(j)}\right)^\top + \frac{1}{n} \sum_{k=1}^{j-1} \left(1 - x^{(k)}\right) \left(1 - x^{(k)}\right)^\top,
\end{aligned}$$

where  $\alpha_j$  is initialized as  $\alpha_{d+1} = \mu_{\min}(M_d)/2$  and is doubled as necessary on each  $j$  iteration to ensure  $x^{(j)} \in \mathcal{X}$ . For specific details regarding our implementation, we refer interested readers to our Python code at <https://anonymous.4open.science/r/SparseContextualCDFRegression-E57B/>.

In our experiments, we compare the  $\ell^2$ -norm estimation error of our proposed lasso estimator (1) against the ridge regression baseline introduced in Zhang et al. (2024, Section 3.1). For both models, we use the regularization hyperparameter  $\lambda = 4\sqrt{(2/n)\log(2d/\delta)}$  specified in Theorem 1, with  $\delta = 0.001$ . Different values of  $\lambda$  and  $\delta$  produced qualitatively similar results. We investigate how the estimation errors scale with various problem dimensions, and report means and standard deviations over 30 independent random trials in Figure 1 for each configuration under consideration. We train both models on the same set of generated samples in each random trial.

Firstly, we investigate the effect of the sample size  $n$ . We choose 100 logarithmically spaced points for  $n$  from  $10^4$  to  $10^6$ , and fix the CDF basis dimension  $d = 10$  and parameter sparsity  $s = 5$ . Figure 1(a) graphs the estimation errors against  $n$  on a log-log plot. The lasso trend line has slope  $-1/2$ , matching the theoretical  $O(1/\sqrt{n})$  bound in Theorem 1 and substantially outperforming the ridge baseline. As a further point of comparison, we repeat this experiment using handpicked regularization hyperparameters for the ridge estimator and plot the results in Figure 1(b). Values of  $\lambda$  less than  $10^{-3}$  or greater than  $10^{-1}$  produced results comparable to the yellow and purple trend lines, respectively. Observe that the accuracy of the ridge estimator with  $\lambda = 10^{-3}$  matches the lasso estimator when  $n \leq 10^5$ , but quickly plateaus for larger sample sizes. This control experiment confirms that the superior accuracy of our lasso estimator in the large  $n$  setting cannot be emulated by the ridge baseline regardless of hyperparameter tuning.

Secondly, we investigate how the estimation errors scale with the sparsity  $s$  of the true parameter  $\theta_*$ . We fix the sample size  $n = 10^5$  and CDF basis dimension  $d = 30$ , and consider  $s$  from 1 to 30. Figure 1(c) graphs the estimation errors against  $s$  on a linear plot. The lasso trend line reflects the theoretical  $O(\sqrt{s})$  bound in Theorem 1, and our lasso estimator notably outperforms the ridge baseline in the sparse regime  $s \ll d$ . To interpret the decreasing ridge trend line, note that  $\|\theta_*\|_1 = 1$  for any value of  $s$  (because  $\theta_*$  is a PMF), and so denser parameter vectors (with greater  $s$ ) tend to have smaller  $\|\theta_*\|$ . Thus, the ridge regularization penalty  $\lambda\|\theta_*\|^2$  at the true parameter vector decreases as  $s$  increases, leading to improved accuracy.

Lastly, we investigate the dependence of the estimation errors on the CDF basis dimension  $d$ . We fix the sample size  $n = 10^5$  and parameter sparsity  $s = 10$ , and consider  $d$  from 10 to 100. Figure 1(d) graphs the estimation errors against  $d$  on a linear plot. The ridge and lasso trend lines indicate the respective theoretical bounds of  $O(\sqrt{d})$  (Zhang et al., 2024, Section 3.2) and  $O(\sqrt{\log d})$  (Theorem 1).

## 8 Conclusion

In this paper, we introduced the task of sparse contextual CDF regression and proposed three basis selection techniques for this problem stemming from the canonical lasso, elastic net, and Dantzig selector regression methods. We derived upper bounds of  $\tilde{O}(\sqrt{s/n})$  and  $\tilde{O}(\sqrt{s}/\sqrt[4]{n})$  on estimation error, and obtained a matching

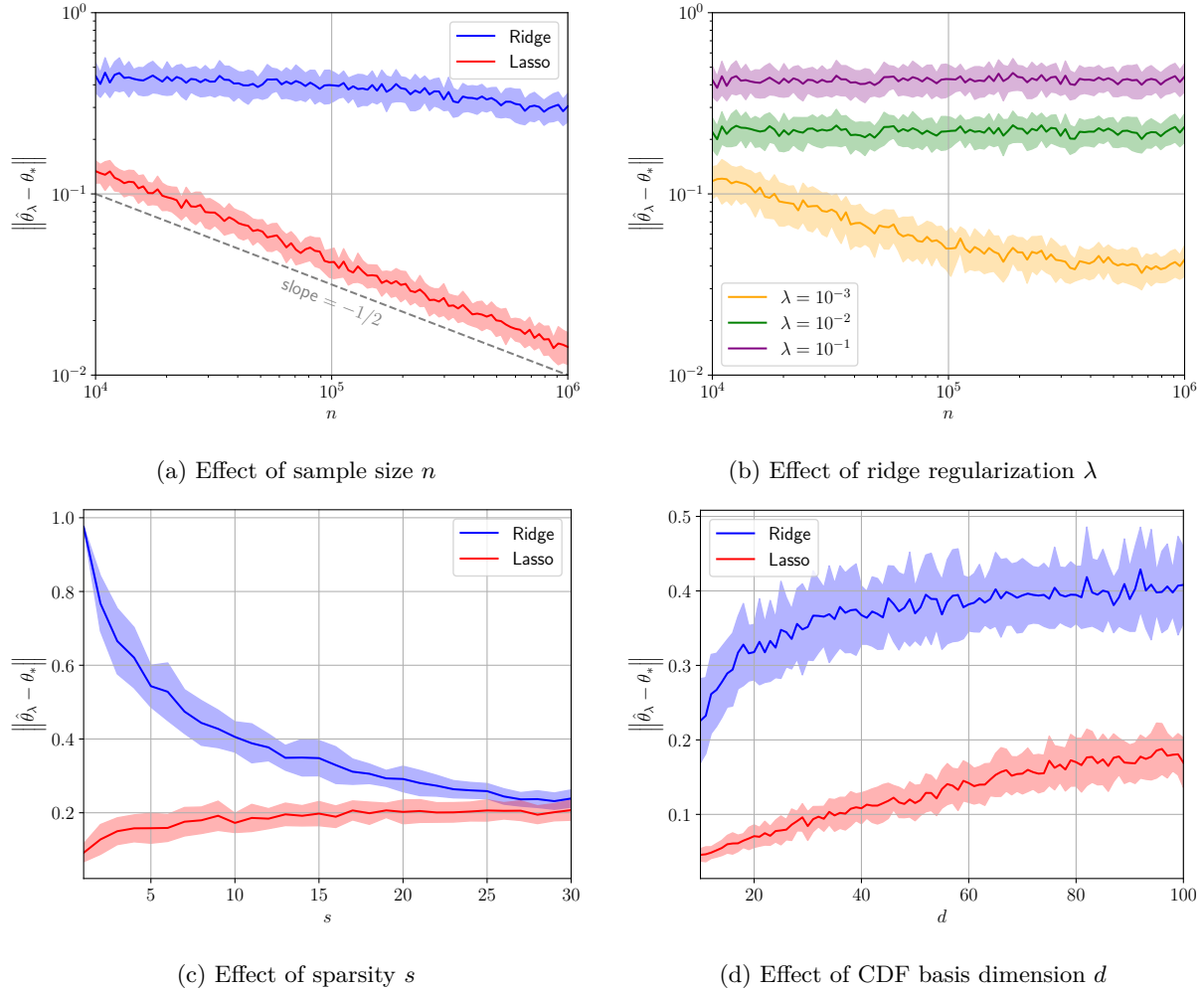


Figure 1: Means and standard deviations of ridge and lasso estimation errors against various hyperparameters in the synthetic Bernoulli experiments. Figures 1(a), 1(c) and 1(d) contrast the ridge and lasso estimators with common regularization hyperparameter  $\lambda = 4\sqrt{(2/n)\log(2d/\delta)}$ , as defined in Theorem 1. For further comparison, Figure 1(b) visualizes the ridge estimation error for various values of  $\lambda$ .

lower bound on minimax risk to establish the optimality of our proposed lasso and Dantzig selector estimators. In particular, our estimation bounds have sub-polynomial dependence on the dimension  $d$  of the CDF regression basis, enabling our methods to perform basis selection with exponentially many irrelevant features and furthering progress towards the ultimate goal of general contextual CDF estimation.

We suggest two directions for future work. Firstly, our present analysis holds only when  $d$  is finite. A natural continuation of our research may investigate similar basis selection methods for CDF regression with infinite-dimensional feature maps. Another promising follow-up direction is to generalize our results to CDF regression with the least absolute deviation ( $\ell^1$ -) loss, in the vein of previous work which combines robust regression and variable selection (Wang et al., 2007).

Overall, our main contributions and proposed future directions indicate that contextual CDF estimation remains a fruitful area of theoretical investigation, accompanied by immediate relevance to a profusion of downstream scientific applications.

## References

- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. *Advances in neural information processing systems*, 24, 2011.
- Carlo Acerbi. Spectral measures of risk: A coherent representation of subjective risk aversion. *Journal of Banking & Finance*, 26(7):1505–1518, 2002.
- Philippe Artzner, Freddy Delbaen, Jean-Marc Eber, and David Heath. Coherent measures of risk. *Mathematical finance*, 9(3):203–228, 1999.
- Afonso S Bandeira, Matthew Fickus, Dustin G Mixon, and Percy Wong. The road to deterministic matrices with the restricted isometry property. *Journal of Fourier Analysis and Applications*, 19(6):1123–1149, 2013.
- Rina Foygel Barber, Matthew Reimherr, and Thomas Schill. The function-on-scalar LASSO with applications to longitudinal GWAS. *Electronic Journal of Statistics*, 11(1):1351 – 1389, 2017. doi: 10.1214/17-EJS1260. URL <https://doi.org/10.1214/17-EJS1260>.
- Amir Beck. *Introduction to nonlinear optimization: Theory, algorithms, and applications with MATLAB*. SIAM, 2014.
- E.J. Candes and T. Tao. Decoding by linear programming. *IEEE Transactions on Information Theory*, 51(12):4203–4215, 2005. doi: 10.1109/TIT.2005.858979.
- Emmanuel Candes and Terence Tao. The Dantzig selector: Statistical estimation when  $p$  is much larger than  $n$ . *The Annals of Statistics*, 35(6):2313 – 2351, 2007. doi: 10.1214/009053606000001523. URL <https://doi.org/10.1214/009053606000001523>.
- Francesco Paolo Cantelli. Sulla determinazione empirica delle leggi di probabilita. *Giorn. Ist. Ital. Attuari*, 4(421-424), 1933.
- Asaf Cassel, Shie Mannor, and Assaf Zeevi. A general framework for bandit problems beyond cumulative objectives. *Mathematics of Operations Research*, 2023.
- Fabio Centofanti, Matteo Fontana, Antonio Lepore, and Simone Vantini. Smooth lasso estimator for the function-on-function linear regression model. *Computational Statistics & Data Analysis*, 176:107556, 2022.
- Luc Devroye, László Györfi, and Gábor Lugosi. *A probabilistic theory of pattern recognition*, volume 31. Springer Science & Business Media, 2013.
- Aryeh Dvoretzky, Jack Kiefer, and Jacob Wolfowitz. Asymptotic minimax character of the sample distribution function and of the classical multinomial estimator. *The Annals of Mathematical Statistics*, pp. 642–669, 1956.
- Valery Glivenko. Sulla determinazione empirica delle leggi di probabilita. *Gion. Ist. Ital. Attauri.*, 4:92–99, 1933.
- Arthur E Hoerl and Robert W Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, 1970.
- Zhaoping Hong and Heng Lian. Inference of genetic networks from time course expression data using functional regression with lasso penalty. *Communications in Statistics-Theory and Methods*, 40(10): 1768–1779, 2011.
- Daniel Hsu, Sham M Kakade, and Tong Zhang. Random design analysis of ridge regression. In *Conference on learning theory*, pp. 9–1. JMLR Workshop and Conference Proceedings, 2012.
- Audrey Huang, Leqi Liu, Zachary Lipton, and Kamyar Azizzadenesheli. Off-policy risk assessment in contextual bandits. *Advances in Neural Information Processing Systems*, 34, 2021.

- Audrey Huang, Leqi Liu, Zachary C Lipton, and Kamyar Azizzadenesheli. Off-policy risk assessment for markov decision processes. In *Artificial Intelligence and Statistics*, 2022.
- Cheng Jie, LA Prashanth, Michael Fu, Steve Marcus, and Csaba Szepesvári. Stochastic optimization in a cumulative prospect theory framework. *IEEE Transactions on Automatic Control*, 63(9):2867–2882, 2018.
- PAVLO A. Krokmal. Higher moment coherent risk measures. *Quantitative Finance*, 7(4):373–387, 2007. doi: 10.1080/14697680701458307.
- Leqi Liu, Audrey Huang, Zachary Lipton, and Kamyar Azizzadenesheli. Supervised learning with general risk functionals. In *International Conference on Machine Learning*, pp. 12570–12592. PMLR, 2022.
- Paolo Maranzano, Philipp Otto, and Alessandro Fassò. Adaptive lasso estimation for functional hidden dynamic geostatistical models. *Stochastic Environmental Research and Risk Assessment*, pp. 1–23, 2023.
- Pascal Massart. The tight constant in the dvoretzky-kiefer-wolfowitz inequality. *The annals of Probability*, pp. 1269–1283, 1990.
- Kevin P Murphy. *Machine learning: a probabilistic perspective*. MIT press, 2012.
- Jasdeep Pannu and Nedret Billor. Robust group-lasso for functional regression model. *Communications in statistics-simulation and computation*, 46(5):3356–3374, 2017.
- LA Prashanth, Cheng Jie, Michael Fu, Steve Marcus, and Csaba Szepesvári. Cumulative prospect theory meets reinforcement learning: Prediction and control. In *International Conference on Machine Learning*, pp. 1406–1415. PMLR, 2016.
- R Tyrrell Rockafellar, Stanislav Uryasev, et al. Optimization of conditional value-at-risk. *Journal of risk*, 2: 21–42, 2000.
- Juan Romo, Rosa E. Lillo, and Nicola Mingotti. Lasso variable selection in functional regression. In *DES - Working Papers. Statistics and Econometrics*, 2013. URL <https://api.semanticscholar.org/CorpusID:118891754>.
- Amir Sani, Alessandro Lazaric, and Rémi Munos. Risk-aversion in multi-armed bandits. *arXiv preprint arXiv:1301.1936*, 2013.
- Alexander Shapiro, Darinka Dentcheva, and Andrzej Ruszczyński. *Lectures on stochastic programming: modeling and theory*. SIAM, 2014.
- Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 58(1):267–288, 1996.
- Joel A Tropp. User-friendly tail bounds for sums of random matrices. *Foundations of computational mathematics*, 12(4):389–434, 2012.
- Sattar Vakili and Qing Zhao. Mean-variance and value at risk in multi-armed bandit problems. In *2015 53rd Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, pp. 1330–1335. IEEE, 2015.
- Martin J Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge university press, 2019.
- Hansheng Wang, Guodong Li, and Guohua Jiang. Robust regression shrinkage and consistent variable selection through the lad-lasso. *Journal of Business & Economic Statistics*, 25(3):347–355, 2007.
- Julia L Wirch and Mary R Hardy. Distortion risk measures: Coherence and stochastic dominance. In *International congress on insurance: Mathematics and economics*, pp. 15–17, 2001.



- William Wong, Audrey Huang, Liu Leqi, Kamyar Azizzadenesheli, and Zachary C Lipton. Riskyzoo: A library for risk-sensitive supervised learning. In *ICML 2022 Workshop on Responsible Decision Making in Dynamic Environments*, 2022.
- Qian Zhang, Anuran Makur, and Kamyar Azizzadenesheli. Functional linear regression of cumulative distribution functions. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=Z0qJCP4eMk>.
- Yihong Zhao, R Todd Ogden, and Philip T Reiss. Wavelet-based lasso in functional linear regression. *Journal of computational and graphical statistics*, 21(3):600–617, 2012.
- Yihong Zhao, Huaihou Chen, and R Todd Ogden. Wavelet-based weighted lasso and screening approaches in functional linear regression. *Journal of Computational and Graphical Statistics*, 24(3):655–675, 2015.
- Alexander Zimin, Rasmus Ibsen-Jensen, and Krishnendu Chatterjee. Generalized risk-aversion in stochastic multi-armed bandits. *arXiv preprint arXiv:1405.0833*, 2014.
- Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 67(2):301–320, 2005.

## A Notation

Let  $\mathbb{N}$  denote the natural numbers starting from 1. Let  $[n] = \mathbb{Z} \cap [1, n]$  denote the set of natural numbers from 1 to  $n$ . Let  $\mathbb{1}\{\cdot\}$  denote the Iverson bracket. When discussing empirical CDFs, we use  $\mathbb{I}_a(t) = \mathbb{1}\{t \geq a\}$  to denote translated unit step functions. Let  $I_d$  denote the  $d \times d$  identity matrix. Let  $\Delta^{d-1} \subset \mathbb{R}^d$  denote the  $(d-1)$ -dimensional probability simplex. Let  $\arg_k \max_{x \in S} f(x)$  denote the top- $k$  maximizers of a function  $f$  over a set  $S$ . In the context of Landau notation, let  $\tilde{O}(\cdot)$  denote asymptotic upper bounds with hidden logarithmic or sub-logarithmic factors. Throughout this paper, all logarithms have base  $e$  and all vectors are column vectors. When invoking a theorem, we use the notation  $x := y$  to pass value  $y$  to variable  $x$  in the theorem statement.

Given a vector  $x \in \mathbb{R}^d$ , let  $[x]_i$  denote its  $i$ th entry. For any set  $S \subseteq [d]$ , let  $[x]_S \in \mathbb{R}^{|S|}$  denote the entries of  $x$  corresponding to the indices in  $S$ . In this context, let  $S^c = [d] - S$  denote set complement with respect to the universe of indices. Let  $\text{supp}(x)$  denote the support of  $x$  and let  $\|x\|_0 = |\text{supp}(x)|$  denote the number of non-zero entries in  $x$ . Let  $\|x\|_p$  denote the  $\ell^p$ -norm of  $x$  for any  $p \in [1, \infty]$ , let  $\|x\|$  denote its Euclidean ( $\ell^2$ -) norm, and let  $\|x\|_A = \sqrt{x^\top A x}$  denote its weighted  $\ell^2$ -norm induced by a positive definite matrix  $A \in \mathbb{R}^{d \times d}$ .

Given a matrix  $A \in \mathbb{R}^{m \times n}$ , let  $[A]_{i,j}$  denote its entry at row  $i$  and column  $j$ . For any set  $S \subseteq [m]$ , let  $[A]_{\langle S \rangle} \in \mathbb{R}^{|S| \times n}$  denote the submatrix obtained from  $A$  by extracting the rows corresponding to the indices in  $S$ . For any sets  $S \subseteq [m]$  and  $T \subseteq [n]$ , let  $[A]_{S,T} \in \mathbb{R}^{|S| \times |T|}$  denote the submatrix obtained from  $A$  by extracting the rows corresponding to the indices in  $S$  and the columns corresponding to the indices in  $T$ . Let  $\|A\|$  denote the induced  $\ell^2$ -operator norm of  $A$ , let  $\mu_{\min}(A)$  denote its minimum eigenvalue, let  $\mu_{\max}(A)$  denote its maximum eigenvalue, and let  $\sigma_{\min}(A)$  denote its minimum singular value.

We denote random variables with uppercase letters and realizations with lowercase letters. We use the notation  $\mathbb{E}_{X \sim \mathcal{P}}[\cdot]$  to emphasize the random variable and distribution which an expectation is taken over. Given probability distributions  $P_1, \dots, P_n$  over random variables  $X_1, \dots, X_n$  respectively, let  $\bigotimes_{j=1}^n P_j$  be the product distribution over  $(X_1, \dots, X_n)$ . Given a probability measure  $\mathbf{m}$  on a set  $S$ , define an inner product between functions  $f, g : S \rightarrow \mathbb{R}$  as  $\langle f, g \rangle = \int_S f(t) g(t) d\mathbf{m}$ , and denote the norm induced by this inner product as  $\|f\|_{\mathcal{L}^2(S, \mathbf{m})} = \sqrt{\langle f, f \rangle}$ . Throughout this paper, integrals and inner products are defined entry-wise for vector-valued and matrix-valued functions.

## B Proof of Lasso Random Design Upper Bound

In this section, we prove Theorem 2.