

# Subquantile Minimization for Kernel Learning in the Huber $\epsilon$ -Contamination Model\*

Arvind Rathnashyam  
RPI Math and CS, [rathna@rpi.edu](mailto:rathna@rpi.edu)

Alex Gittens  
RPI CS, [gittaa@rpi.edu](mailto:gittaa@rpi.edu)

## Abstract

In this paper we propose Subquantile Minimization for learning with adversarial corruption in the training set. Superquantile objectives have been formed in the past in the context of fairness where one wants to learn an underrepresented distribution equally [LPMH21, RRM14]. Our intuition is to learn a more favorable representation of the *majority* class, thus we propose to optimize over the  $p$ -subquantile of the loss in the dataset. In particular, we study the Huber- $\epsilon$  Contamination Problem for Kernel Learning where the distribution is formed as  $\hat{\mathcal{P}} = (1 - \epsilon)\mathcal{P} + \epsilon\mathcal{Q}$ , and we want to find the function  $\inf_{f_{\mathbf{w}} \in \mathcal{H}} \mathbb{E}_{\mathcal{D} \sim \mathcal{P}} [\ell(f_{\mathbf{w}}; \mathbf{X}, \mathbf{y})]$ , from the noisy distribution,  $\hat{\mathcal{P}}$ . We assume the adversary has knowledge of the true distribution of  $\mathcal{P}$ , and is able to corrupt the covariates and the labels of  $\epsilon$  samples. To our knowledge, we are the first to study the problem of general kernel learning in the Huber Contamination Model. In our theoretical analysis, we analyze our non-convex concave objective function with the Moreau Envelope. We show (i) a stationary point with respect to the Moreau Envelope is a good point and (ii) we can reach a stationary point with gradient descent methods. We empirically test Kernel Regression and Kernel Classification on various state of the art datasets and show Subquantile Minimization gives strong results in comparison to the state of the art robust algorithms.

---

\*Preliminary Work

# 1 Introduction

There has been extensive study of algorithms to learn the target distribution from a Huber  $\epsilon$ -Contaminated Model for a Generalized Linear Model (GLM), [DKK<sup>+</sup>19, ADKS22, LBSS21, OZS20, FB81] as well as for linear regression [BJKK17, MGJK19]. Robust Statistics has been studied extensively [DK23] for problems such as high-dimensional mean estimation [PBR19, CDGS20] and Robust Covariance Estimation [CDGW19, FWZ18]. Recently, there has been an interest in solving robust machine learning problems by gradient descent [PSBR18, DKK<sup>+</sup>19]. Subquantile minimization aims to address the shortcomings of standard ERM in applications of noisy/corrupted data [KLA18, JZL<sup>+</sup>18]. In many real-world applications, the covariates have a non-linear dependence on labels [AMMIL12, Section 3.4]. In which case it is suitable to transform the covariates to a different space utilizing kernels [HSS08]. Therefore, in this paper we consider the problem of Robust Learning for Kernel Learning.

**Definition 1** (Huber  $\epsilon$ -Contamination Model [HR09]). Given a corruption parameter  $0 < \epsilon < 0.5$ , a data matrix,  $\mathbf{X}$  and labels  $\mathbf{y}$ . An adversary is allowed to inspect all samples and modify  $\epsilon n$  samples arbitrarily. The algorithm is then given the  $\epsilon$ -corrupted data matrix  $\mathbf{X}$  and  $\mathbf{y}$  as training data.

Current approaches for robust learning across various machine learning tasks often use gradient descent over a robust objective, [LBSS21]. These robust objectives tend to not be convex and therefore do not have a strong analysis on the error bounds for general classes of models.

We similarly propose a robust objective which has a nonconvex-concave objective. This objective has also been proposed recently in [HYwL20] where there has been an analysis in the Binary Classification Task. We show Subquantile Minimization reduces to the same objective in [HYwL20]. We use theory from the weakly-convex concave optimization literature for our error bounds. We are able to leverage this theory by analyzing the asymptotic distribution of a softplus approximation of the Subquantile objective.

The study of Kernel Learning in the Gaussian Design is quite popular, [CLKZ21, Dic16]. In [CLKZ21], the feature space,  $\phi(\mathbf{x}_i) \sim \mathcal{N}(0, \Sigma)$  where  $\Sigma$  is a diagonal matrix of dimension  $p$ , where  $p$  can be infinite. In this work, we adopt a similar framework, and with the power of Mercer’s Theorem [Mer09], we are able to say  $\text{Tr}(\Sigma) < \infty$ . We use this fact extensively in our infinite-dimensional concentration inequalities.

**Theorem 2.** (Informal). Let the dataset be given as  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  such that the labels and features of  $\epsilon n$  samples are arbitrarily corrupted by an adversary.

*Kernelized Regression:*

$$\|\hat{f} - f^*\|_{\mathcal{H}} \leq \varepsilon + O(\sigma) \quad (1)$$

*Kernel Binary Classification:*

$$\|\hat{f} - f^*\|_{\mathcal{H}} \leq \varepsilon + O(\mathcal{E}_{\text{OPT}}) \quad (2)$$

*Kernel Multi-Class Classification:*

$$\|f - f^*\| \leq O(\Xi) \quad (3)$$

## 1.1 Notation

**Reproducing Kernel Hilbert Spaces.** Let the function  $\phi : \mathbb{R}^d \rightarrow \mathcal{H}$  represent the Hilbert Space Representation or ‘feature transform’. We define  $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  as  $k(\mathbf{x}, \mathbf{x}) \triangleq \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_i) \rangle_{\mathcal{H}}$ . For a function in a RKHS,  $f \in \mathcal{H}$ , it follows for some  $\mathbf{w} \in \mathbb{R}^n$ , that  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is defined as  $f(\cdot) \triangleq \sum_{i \in [n]} w_i k(\mathbf{x}_i, \cdot)$ .

## 1.2 Related Work

The idea of iterative thresholding algorithms for robust learning tasks dates back to 1806 by Legendre [Leg06]. From the popularity of Machine Learning, numerous algorithms have been developed in this ideology. Therefore, we will dedicate this section to reviewing such works and to make clear our contributions to the iterative thresholding literature.

Robust Regression via Hard Thresholding [BJK15]. Bhatia et al. study iterative thresholding for least squares regression / sparse recovery. Their theoretical results for the standard gradient descent case cover for known covariance with no feature covariance or Gaussian Noise.

Learning with bad training data via iterative trimmed loss minimization [SS19]. This work considers optimizing over the bottom- $k$  errors by choosing the  $\alpha n$  points with smallest error and then updating the model from these  $\alpha n$ . This general model is the same as ours. Theoretically, this work considers only general linear models.

Trimmed Maximum Likelihood Estimation for Robust Generalized Linear Model [ADKS22]. This work studies a different class of generalized linear models. Interestingly, they show for Gaussian Regression the iterative trimmed maximum likelihood estimator is able to achieve near minimax optimal error. This work does not consider feature corruption and primarily focuses on the covariates sampled with Gaussian Design from Identity covariance.

### 1.3 Contributions

We will now state our main contributions clearly.

1. We provide a novel theoretical framework using the Moreau Envelope for analyzing the iterative trimmed estimator for machine learning tasks.
2. We provide rigorous error bounds for subquantile minimization in the kernel regression, kernel binary classification, and kernel multi-class classification. Furthermore, we provide our bounds for both label and feature corruption with a general Gaussian Design.
3. We perform experiments on state-of-the-art matrices and show the effectiveness of our algorithm compared to other robust learning procedures. Furthermore, we use our experiments to demonstrate the practicality of our theory.

## 2 Subquantile Minimization

We propose to optimize over the subquantile of the risk. The  $p$ -quantile of a random variable,  $U$ , is given as  $\mathcal{Q}_p(U)$ , this is the largest number,  $t$ , such that the probability of  $U \leq t$  is at least  $p$ .

$$\mathcal{Q}_p(U) \leq t \iff \mathbb{P}\{U \leq t\} \geq p \quad (4)$$

The  $p$ -subquantile of the risk is then given by

$$\mathbb{L}_p(U) = \frac{1}{p} \int_0^p \mathcal{Q}_q(U) dq = \mathbb{E}[U | U \leq \mathcal{Q}_p(U)] = \max_{t \in \mathbb{R}} \left\{ t - \frac{1}{p} \mathbb{E}(t - U)^+ \right\} \quad (5)$$

Given an objective function,  $\ell$ , the kernelized learning problem becomes:

$$\min_{f_{\mathbf{w}} \in \mathcal{K}} \max_{t \in \mathbb{R}} \left\{ g(t, f_{\mathbf{w}}) \triangleq t - \sum_{i=1}^n (t - (f_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2)^+ \right\} \quad (6)$$

where  $t$  is the  $p$ -quantile of the empirical risk. Note that for a fixed  $t$  therefore the objective is not concave with respect to  $\mathbf{w}$ . Thus, to solve this problem we use the iterations from Equation 11 in [RHL<sup>+</sup>20]. Let  $\text{Proj}_{\mathcal{K}}$  be the projection of a function on to the convex set  $\mathcal{K} \triangleq \{f \in \mathcal{H} : \|f\|_{\mathcal{H}} \leq R\}$ , then our update steps are

$$t^{(k+1)} = \arg \max_{t \in \mathbb{R}} g(f_{\mathbf{w}}^{(k)}, t) \quad (7)$$

$$f_{\mathbf{w}}^{(k+1)} = \text{Proj}_{\mathcal{K}} \left( f_{\mathbf{w}}^{(k)} - \alpha \nabla_f g(f_{\mathbf{w}}^{(k)}, t^{(k+1)}) \right) \quad (8)$$

We provide an algorithm for Subquantile Minimization of the ridge regression and classification kernel learning algorithm.

### 3 Theory

To consider theoretical guarantees of Subquantile Minimization, we first analyze the inner and outer optimization problems. We first analyze kernel learning in the presence of corrupted data. Next, we provide error bounds for the two most important kernel learning problems, kernel ridge regression, and kernel classification. Now we will give our first result regarding kernel learning in the Huber  $\epsilon$ -contamination model. Now we will analyze the two-step minimax optimization steps described in Equations (7) and (8).

**Lemma 3.** *Let  $f(\mathbf{x}; \mathbf{w})$  be a convex loss function. Let  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  denote the  $n$  data points ordered such that  $f(\mathbf{x}_1; \mathbf{w}, y_1) \leq f(\mathbf{x}_2; \mathbf{w}, y_2) \leq \dots \leq f(\mathbf{x}_n; \mathbf{w}, y_n)$ . If we denote  $\hat{v}_i \triangleq f(\mathbf{x}_i; \mathbf{w}, y_i)$ , it then follows  $\hat{v}_{np} \in \arg \max_{t \in \mathbb{R}} g(t, \mathbf{w})$ .*

Proof is given in Appendix B.1. From Lemma 3, we see that  $t$  will be greater than or equal to the errors of exactly  $np$  points. Thus, we are continuously updating over the  $np$  minimum errors.

**Lemma 4.** *Let  $\hat{v}_i \triangleq f(\mathbf{x}_i; \mathbf{w}, y_i)$  s.t.  $\hat{v}_{i-1} \leq \hat{v}_i \leq \hat{v}_{i+1}$ , if we choose  $t^{(k+1)} = \hat{v}_{np}$  as by Lemma 3, it then follows  $\nabla_{\mathbf{w}} g(t^{(k)}, f_{\mathbf{w}}^{(k)}) = \frac{1}{np} \sum_{i=1}^{np} \nabla f(\mathbf{x}_i; f_{\mathbf{w}}^{(k)}, y_i)$*

Proof is given in Appendix B.2.

#### 3.1 Kernelized Binary Classification

The Negative Log Likelihood for the the Kernel Classification problem is given by the following equation for a single training pair  $(\mathbf{x}_i, y_i)$

$$\ell(\mathbf{x}_i, y_i; \mathbf{w}) = -y_i \log(\sigma(\mathbf{w}^\top \mathbf{x}_i)) - (1 - y_i) \log(1 - \sigma(\mathbf{w}^\top \mathbf{x}_i)) \quad (9)$$

We will now give our algorithm.

**Algorithm 1** (Subquantile Minimization for Binary Classification).

**Input:** Data Matrix:  $\mathbf{X} \in \mathbb{R}^{n \times d}$ ,  $n \gg d$ ; Labels:  $\mathbf{y} \in \mathbb{R}^n$ , Closed and Convex set  $\mathcal{K} \subset \mathcal{H}$

**Output:** Function in  $\mathcal{H}$ :  $\hat{f}$

1. Set the step-size

$$\eta \geq \Omega \left( \frac{\lambda_{\min}(\mathbf{X}^\top \mathbf{X})}{2 \left\| \sum_{i \in X} \mathbf{x}_i \right\|_2^2} \right)$$

2. Set the number of iterations

$$T = O \left( \log \left( \left( \frac{\lambda_{\max}(\Sigma) \|f^*\|_{\mathcal{H}}}{\sqrt{n}} \right) \frac{1}{\epsilon} \right) \right)$$

3. **for**  $k = 1, 2, \dots, T$  **do**

3. Find the Subquantile denoted as  $S^{(k)}$  as the set of  $(1 - \epsilon)n$  elements with the lowest error with respect to the loss function.

4. Calculate the gradient update.

$$\nabla_f g(t^{(k+1)}, f^{(k)}) \leftarrow \frac{1}{n(1 - \epsilon)} \sum_{i \in S^{(k)}} (\sigma(f^{(k)}(\mathbf{x}_i) - y_i) \cdot \mathbf{x}_i^\top$$

5. Perform Projected Gradient Descent Iteration

$$f^{(k+1)} \leftarrow \text{Proj}_{\mathcal{K}} \left[ f^{(k)} - \eta \nabla g(f^{(k)}, t^{(k+1)}) \right]$$

**Return:** Function in  $\mathcal{H}$ :  $f^{(T)}$

**Theorem 5** (Subquantile Minimization for Binary Classification is Good with High Probability). *Let Algorithm 1 be run on a dataset  $\mathcal{D} \sim \hat{\mathcal{P}}$  with learning rate  $\eta \triangleq \Omega(L^{-1})$ . Then after  $O\left(\log\left(\left(\frac{\lambda_{\max}(\mathbf{\Sigma})\|f^*\|_{\mathcal{H}}}{\sqrt{n}}\right)\frac{1}{\epsilon}\right)\right)$  gradient descent iterations,*

$$\|f^{(T)} - f^*\|_2 \leq \epsilon + O\left(\frac{\mathcal{E}_{\text{OPT}}}{\sqrt{n(1-\epsilon)}}\right) \quad (10)$$

where  $\mathcal{E}_{\text{OPT}} \triangleq \sum_{i \in \mathcal{P}} \left( \mathbf{Pr}_{(\mathbf{x}_i, y_i) \sim \mathcal{P}} \{y_i \mid \mathbf{x}_i\} - y_i \right)^2$ .

The full proof is given in Appendix C.2. In Theorem 5, we introduce  $\mathcal{E}_{\text{OPT}}$ , which says we are only able to learn up to the intrinsic noise within the target function.

### 3.2 Kernelized Multi-Class Classification

The Negative Log-Likelihood Loss for the the Kernel Multi-Class Classification problem is given by the following equation for a single training pair  $(\mathbf{x}_i, y_i)$ , note  $\mathbf{W} \in \mathbb{R}^{n \times |\mathcal{Y}|} = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_{|\mathcal{Y}|}]$ .

$$\ell(\mathbf{x}_i, y_i; \mathbf{f}_{\mathbf{W}}) = - \sum_{j=1}^{|\mathcal{Y}|} \mathbb{I}\{j = y_i\} \log \left( \frac{\exp(f_{\mathbf{w}_j}(\mathbf{x}_i))}{\sum_{k=1}^{|\mathcal{Y}|} \exp(f_{\mathbf{w}_k}(\mathbf{x}_i))} \right) \quad (11)$$

We will now give the algorithm.

**Algorithm 2** (Subquantile Minimization for Binary Classification).

**Input:** Data Matrix:  $\mathbf{X} \in \mathbb{R}^{n \times d}$ ,  $n \gg d$ ; Labels:  $\mathbf{y} \in \mathbb{R}^{n \times |\mathcal{Y}|}$ , Closed and Convex set  $\mathcal{K} \subset \mathcal{H}$

**Output:** Functions in  $\mathcal{H}$ :  $\hat{f}_{\ell}$  for  $\ell \in [|\mathcal{Y}|]$

1. Set the step-size

$$\eta \leq O(\Xi)$$

2. Set the number of iterations

$$T = O(Xi)$$

3. **for**  $k = 1, 2, \dots, T$  **do**

3. Find the Subquantile denoted as  $S^{(k)}$  as the set of  $(1 - \epsilon)n$  elements with the lowest error with respect to the loss function.

4. Calculate the gradient update for each function.

$$\nabla_{f_{g_{\ell}}(t^{(k+1)}, f^{(k)})} \leftarrow \frac{1}{n(1-\epsilon)} \sum_{i \in S^{(k)}} (\text{softmax}(f_{\ell}^{(k)}(\mathbf{x}_i) - y_i) \cdot \mathbf{x}_i^{\top}$$

5. Perform Projected Gradient Descent Iteration for each function.

$$f_{\ell}^{(k+1)} \leftarrow \text{Proj}_{\mathcal{K}} \left[ f_{\ell}^{(k)} - \eta \nabla_{g_{\ell}}(f_{\ell}^{(k)}, t^{(k+1)}) \right]$$

**Return:** Function in  $\mathcal{H}$ :  $f^{(T)}$

**Theorem 6** (Subquantile Minimization for Kernelized Multi-Class Classification is Good with High Probability). *Let Algorithm 1 be run on a dataset  $\mathcal{D} \sim \hat{\mathcal{P}}$  with learning rate  $\eta \triangleq \beta^{-1}$ . Then,*

$$\|f_{\mathbf{W}} - f_{\mathbf{W}}^*\|_{\mathcal{H}} \leq O(\Xi) \quad (12)$$

## 4 Discussion

The main contribution of this paper is the study of a nonconvex-concave formulation of Subquantile minimization for the robust learning problem for kernel ridge regression and kernel classification. We present an algorithm to solve the nonconvex-concave formulation and prove rigorous error bounds which show that the more good data that is given decreases the error bounds. We also present accelerated gradient methods for the two-step algorithm to solve the nonconvex-concave optimization problem and give novel theoretical bounds.

**Theory.** We develop strong theoretical bounds on the normed difference between the function returned by Subquantile Minimization and the optimal function for data in the target distribution,  $\mathbb{P}$ , in the Gaussian Design. In expectation and with high probability, given sufficient data dependent on the kernel, we obtain a near minimax optimal error bound for a general positive definite continuous kernel. Our theoretical analysis is novel in that it utilizes the Moreau Envelope from a min-max formulation of the iterative thresholding algorithm.

**Experiments.** From our experiments, we see Subquantile Minimization is competitive with algorithms developed solely for robust linear regression as well as other meta-algorithms. Our theoretical analysis is through the lens of kernel-learning, but the generalization to linear regression from a non-kernel perspective can be done. In kernelized regression, we see SUBQUANTILE is the strongest of the meta-algorithms. Furthermore, in binary and multi-class classification, SUBQUANTILE is very strong. Thus, we can see empirically SUBQUANTILE is the strongest meta-algorithm across all kernelized regression and classification tasks and also the strongest algorithm in linear regression.

**Interpretability.** One of the strengths in Subquantile Optimization is the high interpretability. Once training is finished, we can see the  $n(1-p)$  points with highest error to find the outliers and the features follow Gaussian Design. Furthermore, there is only hyperparameter  $p$ , which should be chosen to be approximately the percentage of inliers in the data and thus is not very difficult to tune for practical purposes. Our theory suggests for a problem where the amount of corruptions is unknown,

**General Assumptions.** The general assumption is the majority of the data should inliers. This is not a very strong assumption, as by the definition of outlier it should be in the minority. Furthermore, we assume the feature maps have a Gaussian Design. Such a design in many prior works in kernel learning and we therefore find it suitable.

**Future Work.** The analysis of Subquantile Minimization can be extended to neural networks as kernel learning can be seen as a one-layer network. This generalization will be appear in subsequent work. Another interesting direction work in optimization is for accelerated methods for optimizing non-convex concave min-max problems with a maximization oracle. The current theory analyzes standard gradient descent for the minimization. Ideas such as Momentum and Nesterov Acceleration in conjunction with the maximum oracle are interesting and can be analyzed in future work.

## References

- [ADKS22] Pranjali Awasthi, Abhimanyu Das, Weihao Kong, and Rajat Sen. Trimmed maximum likelihood estimation for robust generalized linear model. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. 1, 2
- [AMMIL12] Yaser S Abu-Mostafa, Malik Magdon-Ismael, and Hsuan-Tien Lin. *Learning from data*, volume 4. AMLBook New York, 2012. 1
- [BJK15] Kush Bhatia, Prateek Jain, and Purushottam Kar. Robust regression via hard thresholding. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015. 1
- [BJKK17] Kush Bhatia, Prateek Jain, Parameswaran Kamalaruban, and Purushottam Kar. Consistent robust regression. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. 1

- [CDGS20] Yu Cheng, Ilias Diakonikolas, Rong Ge, and Mahdi Soltanolkotabi. High-dimensional robust mean estimation via gradient descent. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 1768–1778. PMLR, 13–18 Jul 2020. [1](#)
- [CDGW19] Yu Cheng, Ilias Diakonikolas, Rong Ge, and David P. Woodruff. Faster algorithms for high-dimensional robust covariance estimation. In Alina Beygelzimer and Daniel Hsu, editors, *Proceedings of the Thirty-Second Conference on Learning Theory*, volume 99 of *Proceedings of Machine Learning Research*, pages 727–757. PMLR, 25–28 Jun 2019. [1](#)
- [CLKZ21] Hugo Cui, Bruno Loureiro, Florent Krzakala, and Lenka Zdeborová. Generalization error rates in kernel regression: The crossover from the noiseless to noisy regime. *Advances in Neural Information Processing Systems*, 34:10131–10143, 2021. [1](#)
- [Dic16] Lee H Dicker. Ridge regression and asymptotic minimax estimation over spheres of growing dimension. 2016. [1](#)
- [DK23] Ilias Diakonikolas and Daniel M Kane. *Algorithmic high-dimensional robust statistics*. Cambridge University Press, 2023. [1](#)
- [DKK<sup>+</sup>19] Ilias Diakonikolas, Gautam Kamath, Daniel M. Kane, Jerry Li, Jacob Steinhardt, and Alistair Stewart. Sever: A robust meta-algorithm for stochastic optimization. In *Proceedings of the 36th International Conference on Machine Learning, ICML '19*, pages 1596–1606. JMLR, Inc., 2019. [1](#)
- [FB81] Martin A. Fischler and Robert C. Bolles. Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography. *Commun. ACM*, 24(6):381–395, jun 1981. [1](#)
- [FWZ18] Jianqing Fan, Weichen Wang, and Yiqiao Zhong. An  $l$  eigenvector perturbation bound and its application to robust covariance estimation. *Journal of Machine Learning Research*, 18(207):1–42, 2018. [1](#)
- [Gre13] Arthur Gretton. Introduction to rkhs, and some simple kernel algorithms. *Adv. Top. Mach. Learn. Lecture Conducted from University College London*, 16(5-3):2, 2013. [9](#)
- [HR09] Peter J. Huber and Elvezio. Ronchetti. *Robust statistics*. Wiley series in probability and statistics. Wiley, Hoboken, N.J., 2nd ed. edition, 2009. [1](#)
- [HSS08] Thomas Hofmann, Bernhard Schölkopf, and Alexander J. Smola. Kernel methods in machine learning. *The Annals of Statistics*, 36(3):1171 – 1220, 2008. [1](#)
- [HYwL20] Shu Hu, Yiming Ying, xin wang, and Siwei Lyu. Learning by minimizing the sum of ranked range. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 21013–21023. Curran Associates, Inc., 2020. [1](#)
- [JZL<sup>+</sup>18] Lu Jiang, Zhengyuan Zhou, Thomas Leung, Li-Jia Li, and Li Fei-Fei. Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels. In *ICML*, 2018. [1](#)
- [KLA18] Ashish Khetan, Zachary C. Lipton, and Anima Anandkumar. Learning from noisy singly-labeled data. In *International Conference on Learning Representations*, 2018. [1](#)
- [LBSS21] Tian Li, Ahmad Beirami, Maziar Sanjabi, and Virginia Smith. Tilted empirical risk minimization. In *International Conference on Learning Representations*, 2021. [1](#)
- [Leg06] Adrien M Legendre. *Nouvelles methodes pour la determination des orbites des cometes: avec un supplement contenant divers perfectionnemens de ces methodes et leur application aux deux cometes de 1805*. Courcier, 1806. [1](#)

- [LPMH21] Yassine Laguel, Krishna Pillutla, Jérôme Malick, and Zaid Harchaoui. Superquantiles at work: Machine learning applications and efficient subgradient computation. *Set-Valued and Variational Analysis*, 29(4):967–996, Dec 2021.
- [Mer09] James Mercer. Xvi. functions of positive and negative type, and their connection the theory of integral equations. *Philosophical transactions of the royal society of London. Series A, containing papers of a mathematical or physical character*, 209(441-458):415–446, 1909. [1](#)
- [MGJK19] Bhaskar Mukhoty, Govind Gopakumar, Prateek Jain, and Purushottam Kar. Globally-convergent iteratively reweighted least squares for robust regression problems. In Kamalika Chaudhuri and Masashi Sugiyama, editors, *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pages 313–322. PMLR, 16–18 Apr 2019. [1](#)
- [OZS20] Muhammad Osama, Dave Zachariah, and Petre Stoica. Robust risk minimization for statistical learning from corrupted data. *IEEE Open Journal of Signal Processing*, 1:287–294, 2020. [1](#)
- [PBR19] Adarsh Prasad, Sivaraman Balakrishnan, and Pradeep Ravikumar. A unified approach to robust mean estimation. *arXiv preprint arXiv:1907.00927*, 2019. [1](#)
- [PSBR18] Adarsh Prasad, Arun Sai Suggala, Sivaraman Balakrishnan, and Pradeep Ravikumar. Robust estimation via robust gradient estimation. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 82, 2018. [1](#)
- [RHL<sup>+</sup>20] Meisam Razaviyayn, Tianjian Huang, Songtao Lu, Maher Nouiehed, Maziar Sanjabi, and Mingyi Hong. Nonconvex min-max optimization: Applications, challenges, and recent theoretical advances. *IEEE Signal Processing Magazine*, 37(5):55–66, 2020. [2](#)
- [RRM14] R.T. Rockafellar, J.O. Royset, and S.I. Miranda. Superquantile regression with applications to buffered reliability, uncertainty quantification, and conditional value-at-risk. *European Journal of Operational Research*, 234(1):140–154, 2014.
- [SS19] Yanyao Shen and Sujay Sanghavi. Learning with bad training data via iterative trimmed loss minimization. In *International Conference on Machine Learning*, pages 5739–5748. PMLR, 2019. [2](#)
- [Wey12] Hermann Weyl. Das asymptotische verteilungsgesetz der eigenwerte linearer partieller differentialgleichungen (mit einer anwendung auf die theorie der hohlraumstrahlung). *Mathematische Annalen*, 71(4):441–479, 1912. [14](#), [16](#)
- [You12] William Henry Young. On classes of summable functions and their fourier series. *Proceedings of the Royal Society of London. Series A, Containing Papers of a Mathematical and Physical Character*, 87(594):225–229, 1912. [13](#)



## A Probability Theory

In this section we will give various concentration inequalities on the inlier data for functions in the Reproducing Kernel Hilbert Space. We will first give our assumptions for robust kernelized regression.

**Assumption 7** (Gaussian Design). We assume for  $\mathbf{x}_i \sim \mathcal{P} \in \mathcal{X}$ , then it follows for the feature map,  $\phi(\cdot) : \mathcal{X} \rightarrow \mathcal{H}$ ,

$$\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}) \quad (13)$$

where  $\mathbf{\Sigma}$  is a possibly infinite dimensional covariance operator.

**Assumption 8** (Normal Residuals). The residual is defined as  $\mu_i \triangleq f_{\mathbf{w}}^*(\mathbf{x}_i) - y_i$ . Then we assume for some  $\sigma > 0$ , it follows

$$\mu_i \sim \mathcal{N}(0, \sigma^2) \quad (14)$$

**Proposition 9** (Expected Maximum  $P_k$ ). Let  $\mathbf{x}_i \sim \mathcal{P}$  such that  $\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})$  (Assumption 7). Then it follows for any  $s \geq 1$

$$\mathbf{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \left[ \max_{i \in [n]} k(\mathbf{x}_i, \mathbf{x}_i) \right] \leq 2 \text{Tr}(\mathbf{\Sigma}) \log \left( \frac{n \cdot s}{2 \text{Tr}(\mathbf{\Sigma})} \right) + \frac{1}{s} \quad (15)$$

**Proof.** We will use the integral identity of the expectation of a random variable to make our claim. Throughout the proof, let  $C$  be a positive to be determined constant.

$$\mathbf{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \left[ \max_{i \in [n]} k(\mathbf{x}_i, \mathbf{x}_i) \right] \leq C + \int_C^\infty \mathbf{Pr}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \left\{ \max_{i \in [n]} k(\mathbf{x}_i, \mathbf{x}_i) \geq t \right\} dt \quad (16)$$

$$\stackrel{(i)}{\leq} C + n \int_C^\infty \mathbf{Pr}_{\phi(\mathbf{x}) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \{k(\mathbf{x}, \mathbf{x}) \geq t\} dt \quad (17)$$

$$\stackrel{(ii)}{=} C + n \int_C^\infty \mathbf{Pr}_{\phi(\mathbf{x}) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \left\{ \|\phi(\mathbf{x})\|_{\mathcal{H}} \geq \sqrt{t} \right\} dt \quad (18)$$

$$\stackrel{(iii)}{=} C + n \int_C^\infty \mathbf{Pr}_{\psi(\mathbf{x}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left\{ \left\| \mathbf{\Psi} \mathbf{\Lambda}^{1/2} \psi(\mathbf{x}) \right\|_{\mathcal{H}} \geq \sqrt{t} \right\} dt \quad (19)$$

$$\leq C + n \int_C^\infty \mathbf{Pr}_{\psi_i(\mathbf{x}) \sim \mathcal{N}(0, 1)} \left\{ \sum_{i=1}^p \sqrt{\lambda_i} \psi_i(\mathbf{x}) \geq \sqrt{t} \right\} dt \quad (20)$$

$$\leq C + n \int_C^\infty \inf_{\theta > 0} \mathbf{E}_{\psi(\mathbf{x}) \sim \mathcal{N}(0, 1)} \left[ \prod_{i=1}^p \exp(\theta \sqrt{\lambda_i} \psi_i(\mathbf{x})) \right] \exp(-\theta \sqrt{t}) dt \quad (21)$$

$$= C + n \int_C^\infty \inf_{\theta > 0} \exp \left[ \theta^2 \frac{\text{Tr}(\mathbf{\Sigma})}{2} - \theta \sqrt{t} \right] dt = C + n \int_C^\infty \exp \left[ -\frac{t}{2 \text{Tr}(\mathbf{\Sigma})} \right] dt \quad (22)$$

$$= C + \frac{n}{2 \text{Tr}(\mathbf{\Sigma})} \exp \left[ -\frac{C}{2 \text{Tr}(\mathbf{\Sigma})} \right] \quad (23)$$

See (i) from the proof of ???. (ii) follows from the reproducing property. In (iii) we define  $\psi(\mathbf{x})$  as the whitened RKHS function. Setting  $C \triangleq 2 \text{Tr}(\mathbf{\Sigma}) \log(s \cdot n / (2 \text{Tr}(\mathbf{\Sigma})))$  completes the proof. ■

**Lemma 10** (Norm of Functions with Gaussian Design in the Reproducing Kernel Hilbert Space). Let  $\mathbf{x}_i \sim \mathbb{P}$  such that  $\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})$  from Assumption 7 and Assumption 8. Then, it follows

$$\mathbf{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \mathbf{E}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \left\| \sum_{i=1}^n \mu_i \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \leq O \left( \sigma \sqrt{n \log n \text{Tr}(\mathbf{\Sigma})} \right) \quad (24)$$

**Proof.** Our proof follows standard ideas from High-Dimensional Probability. Let  $\xi_i$  for  $i \in [n]$  denote i.i.d Rademacher variables such that for  $\xi_i \sim \mathcal{R}$ , it follows  $\mathbb{P}\{\xi_i = 1\} = \mathbb{P}\{\xi_i = -1\} = \frac{1}{2}$ . We then have,

$$\mathbf{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \mathbf{E}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \left\| \sum_{i=1}^n \mu_i \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \leq \mathbf{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma})} \mathbf{E}_{\mu_i \sim \mathcal{N}(0, \sigma^2)} \max_{i \in [n]} |\mu_i| \left\| \sum_{i=1}^n \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \quad (25)$$

$$\stackrel{(i)}{\leq} O\left(\sigma\sqrt{\log n}\right) \mathbf{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbf{E}_{\xi_i \sim \mathcal{R}} \left\| \sum_{i=1}^n \xi_i \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \quad (26)$$

$$\stackrel{(ii)}{\leq} O\left(\sigma\sqrt{\log n}\right) \left( \mathbf{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbf{E}_{\xi_i \sim \mathcal{R}} \left\| \sum_{i=1}^n \xi_i \phi(\mathbf{x}_i) \right\|_{\mathcal{H}}^2 \right)^{1/2} \quad (27)$$

$$= O\left(\sigma\sqrt{\log n}\right) \left( \mathbf{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbf{E}_{\xi_i \sim \mathcal{R}} \left\langle \sum_{i=1}^n \xi_i \phi(\mathbf{x}_i), \sum_{j=1}^n \xi_j \phi(\mathbf{x}_j) \right\rangle_{\mathcal{H}} \right)^{1/2} \quad (28)$$

$$\stackrel{(iii)}{=} O\left(\sigma\sqrt{\log n}\right) \left( \mathbf{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbf{E}_{\xi_i \sim \mathcal{R}} \sum_{i=1}^n \sum_{j=1}^n \xi_i \xi_j k(\mathbf{x}_i, \mathbf{x}_j) \right)^{1/2} \quad (29)$$

$$\stackrel{(iv)}{=} O\left(\sigma\sqrt{\log n}\right) \left( \mathbf{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \sum_{i=1}^n k(x_i, x_i) \right)^{1/2} \quad (30)$$

$$= O\left(\sigma\sqrt{n \log n \operatorname{Tr}(\Sigma)}\right) \quad (31)$$

(i) follows from applying ?? . (ii) follows from Jensen's Inequality. (iii) follows from the definition of the kernel [Gre13]. (iv) holds as we have  $\mathbb{E}[\xi_i \xi_j] = \delta_{i,j}$ , where  $\delta$  is the Kronecker Delta function. ■

**Proposition 11** (Probabilistic bound on Norm of Functions with Gaussian Design in the Reproducing Kernel Hilbert Space). *Let  $\mathbf{x}_i \sim \mathcal{P}$  such that  $\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)$  (Assumption 7). Then it follows*

$$\mathbf{Pr}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\{ \left\| \sum_{i=1}^n \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \geq \sqrt{n \operatorname{Tr}(\Sigma)} \cdot u \right\} \leq e^{-u^2/2} \quad (32)$$

**Proof.** Our proof will utilize a symmetrization argument similar to our previous expected covariance approximation proof, the proof then follows similarly to ?? . on the Let  $C$  be a positive constant to be determined and  $u \geq 1$ . We then have

$$\mathbf{Pr}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\{ \left\| \sum_{i=1}^n \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \geq C \cdot u \right\} = \mathbf{Pr}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\{ \mathbf{E}_{\xi_i \sim \mathcal{R}} \left\| \sum_{i=1}^n \xi_i \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \geq C \cdot u \right\} \quad (33)$$

$$\stackrel{(i)}{\leq} \mathbf{Pr}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\{ \mathbf{E}_{\xi_i \sim \mathcal{R}} \left\| \sum_{i=1}^n \xi_i \phi(\mathbf{x}_i) \right\|_{\mathcal{H}}^2 \geq C^2 \cdot u^2 \right\} \quad (34)$$

$$= \mathbf{Pr}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\{ \mathbf{E}_{\xi_i \sim \mathcal{R}} \sum_{i=1}^n \sum_{j=1}^n \xi_i \xi_j k(\mathbf{x}_i, \mathbf{x}_j) \geq C^2 \cdot u^2 \right\} \quad (35)$$

$$= \mathbf{Pr}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\{ \sum_{i=1}^n k(\mathbf{x}_i, \mathbf{x}_i) \geq C^2 \cdot u^2 \right\} \leq \mathbf{Pr}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\{ \sum_{i=1}^n \|\phi(\mathbf{x}_i)\|_{\mathcal{H}} \geq C \cdot u \right\} \quad (36)$$

$$\stackrel{(ii)}{\leq} \inf_{\theta > 0} \prod_{i=1}^n \prod_{j=1}^p \exp \left[ \frac{\theta^2 \lambda_j}{2} \right] \exp [-\theta C \cdot u] = \exp \left[ -\frac{C^2 \cdot u^2}{2n \operatorname{Tr}(\Sigma)} \right] \quad (37)$$

In (i) we use Jensen's Inequality. For (ii) see the proof for ?? with an additional product term over  $n$ . Finally, setting  $C = \sqrt{n \operatorname{Tr}(\Sigma)}$  completes the proof.

**Proposition 12** (Probabilistic Bound on Infinite Dimensional Covariance Estimation in the Hilbert-Schmidt Norm). *Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  be i.i.d sampled from  $\mathbb{P}$  such that  $\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)$  (Assumption 7), we then have for any  $u \geq 1$*

$$\mathbf{Pr}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\{ \left\| \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) - \Sigma \right\|_{\text{HS}} \geq \frac{\sqrt{3} \operatorname{Tr}(\Sigma)}{\sqrt{n}} \cdot u \right\} \leq e^{-u^2 \operatorname{Tr}(\Sigma)/2} \quad (38)$$

**Proof.** Let  $C$  be a to be determined positive constant and  $u$  be a positive constant.

$$\Pr_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\{ \left\| \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) - \Sigma \right\|_{\text{HS}} \geq C \cdot u \right\} \quad (39)$$

$$\stackrel{(i)}{\leq} \Pr_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\{ \mathbf{E}_{\tilde{\phi}(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\| \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) - \tilde{\phi}(\mathbf{x}_i) \otimes \tilde{\phi}(\mathbf{x}_i) \right\|_{\text{HS}} \geq C \cdot u \right\} \quad (40)$$

$$\stackrel{(ii)}{=} \Pr_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\{ \mathbf{E}_{\tilde{\phi}(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbf{E}_{\xi_i \sim \mathcal{R}} \left\| \frac{1}{n} \sum_{i=1}^n \xi_i (\phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) - \tilde{\phi}(\mathbf{x}_i) \otimes \tilde{\phi}(\mathbf{x}_i)) \right\|_{\text{HS}} \geq C \cdot u \right\} \quad (41)$$

$$\leq \Pr_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\{ \mathbf{E}_{\xi_i \sim \mathcal{R}} \left\| \frac{1}{n} \sum_{i=1}^n \xi_i \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) \right\|_{\text{HS}} \right. \quad (42)$$

$$\left. + \mathbf{E}_{\tilde{\phi}(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbf{E}_{\xi_i \sim \mathcal{R}} \left\| \frac{1}{n} \sum_{i=1}^n \xi_i \tilde{\phi}(\mathbf{x}_i) \otimes \tilde{\phi}(\mathbf{x}_i) \right\|_{\text{HS}} \geq C \cdot u \right\}$$

$$\leq \Pr_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\{ \left( \mathbf{E}_{\xi_i \sim \mathcal{R}} \left\| \frac{1}{n} \sum_{i=1}^n \xi_i \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) \right\|_{\text{HS}}^2 \right)^{1/2} \right. \quad (43)$$

$$\left. + \left( \mathbf{E}_{\tilde{\phi}(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbf{E}_{\xi_i \sim \mathcal{R}} \left\| \frac{1}{n} \sum_{i=1}^n \xi_i \tilde{\phi}(\mathbf{x}_i) \otimes \tilde{\phi}(\mathbf{x}_i) \right\|_{\text{HS}}^2 \right)^{1/2} \geq C \cdot u \right\}$$

We will expand the second term in Equation (43). Let  $e_k$  for  $k \in [p]$  represent an orthonormal basis for the Hilbert Space  $\mathcal{H}$ . By expanding out the Hilbert-Schmidt Norm, we then have

$$\frac{1}{n} \left( \mathbf{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbf{E}_{\xi_i \sim \mathcal{R}} \left\| \sum_{i=1}^n \xi_i \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) \right\|_{\text{HS}}^2 \right)^{1/2}$$

$$= \frac{1}{n} \left( \mathbf{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbf{E}_{\xi_i \sim \mathcal{R}} \sum_{k=1}^p \left\langle \sum_{i=1}^n \xi_i \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) e_k, \sum_{j=1}^n \xi_j \phi(\mathbf{x}_j) \otimes \phi(\mathbf{x}_j) e_k \right\rangle_{\mathcal{H}} \right)^{1/2} \quad (44)$$

$$= \frac{1}{n} \left( \mathbf{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbf{E}_{\xi_i \sim \mathcal{R}} \sum_{k=1}^p \sum_{i=1}^n \sum_{j=1}^n \xi_i \xi_j \langle \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) e_k, \phi(\mathbf{x}_j) \otimes \phi(\mathbf{x}_j) e_k \rangle_{\mathcal{H}} \right)^{1/2} \quad (45)$$

$$\stackrel{(iv)}{\leq} \frac{2}{n} \left( \mathbf{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \sum_{k=1}^p \sum_{i=1}^n \langle \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) e_k, \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) e_k \rangle_{\mathcal{H}} \right)^{1/2} \quad (46)$$

$$= \frac{1}{n} \left( \sum_{i=1}^n \mathbf{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \|\phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i)\|_{\text{HS}}^2 \right)^{1/2} \stackrel{(v)}{=} \frac{2}{\sqrt{n}} \mathbf{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \|\phi(\mathbf{x}_i)\|_{\mathcal{H}}^4 \quad (47)$$

$$= \frac{1}{\sqrt{n}} \mathbf{E}_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} [k^2(\mathbf{x}_i, \mathbf{x}_i)] = \frac{1}{\sqrt{n}} \left( 2 \text{Tr}(\Sigma^2) + \text{Tr}(\Sigma)^2 \right)^{1/2} \leq \sqrt{3} n^{-1/2} \text{Tr}(\Sigma) \quad (48)$$

(iv) follows from noticing  $\mathbf{E}_{\xi_i, \xi_j \sim \mathcal{R}} [\xi_i \xi_j] = \delta_{ij}$ . (v) follows from expanding the Hilbert-Schmidt Norm and applying Parseval's Identity. We note  $\text{Tr}(\Sigma) < \infty$  and therefore even though the covariance operator is infinite-dimensional we are able to get a finite bound on the covariance approximation.

$$(43) \text{ RHS} \leq \Pr_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\{ \frac{\sqrt{\text{Tr}(\mathbf{K})}}{n} + \frac{\sqrt{3} \text{Tr}(\Sigma)}{\sqrt{n}} \geq C \cdot u \right\} \quad (49)$$

$$\leq \Pr_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\{ \sum_{i=1}^n \|\phi(\mathbf{x}_i)\|_{\mathcal{H}} \geq nC \cdot u - \sqrt{3n} \text{Tr}(\Sigma) \right\} \quad (50)$$

$$\leq \exp \left[ - \frac{(nC \cdot u + \sqrt{3n} \text{Tr}(\Sigma))^2}{n \text{Tr}(\Sigma)} \right] \stackrel{(vi)}{\leq} e^{-u^2 \text{Tr}(\Sigma)/2} \quad (51)$$

In (vi) we chose  $C \triangleq \sqrt{3} \text{Tr}(\Sigma)/\sqrt{n}$  and then simplify the resultant probability bound.  $\blacksquare$

**Proposition 13** (Probabilistic Bound on Infinite Dimensional Covariance Estimation in the Hilbert-Schmidt Norm). *Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  be i.i.d sampled from  $\mathcal{P}$  such that  $\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)$  (Assumption 7), we then have for any  $u \geq 1$*

$$\Pr_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\{ \left\| \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) - \Sigma \right\|_{\text{op}} \geq \frac{\|\Sigma\|}{\sqrt{n}} \cdot u \right\} \leq e^{-u^2 \|\Sigma\|/2} \quad (52)$$

**Proof.**<sup>1</sup> Let  $C$  be a to be determined positive constant, and  $u$  be a positive constant.

$$\begin{aligned} \Pr_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\{ \left\| \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) - \Sigma \right\|_{\text{op}} \geq C \cdot u \right\} &\leq \Pr_{\phi(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \left\{ \mathbf{E}_{\xi \sim \mathcal{R}} \left\| \sum_{i=1}^n \xi_i \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) \right\|_{\text{op}} \right. \\ &\quad \left. + \mathbf{E}_{\tilde{\phi}(\mathbf{x}_i) \sim \mathcal{N}(\mathbf{0}, \Sigma)} \mathbf{E}_{\xi_i \sim \mathcal{R}} \left\| \frac{1}{n} \sum_{i=1}^n \xi_i \tilde{\phi}(\mathbf{x}_i) \otimes \tilde{\phi}(\mathbf{x}_i) \right\|_{\text{HS}} \geq C \cdot u \right\} \end{aligned} \quad (53)$$

## B Proofs for Structural Results

In this section we give the deferred proofs of our main structural results of the subquantile objective function.

### B.1 Proof of Lemma 3

**Proof.** First we can note, the max value of  $t$  for  $g$  is equivalent to the min value of  $t$  for  $g$ . We can now find the Fermat Optimality Conditions for  $g$ .

$$\partial(-g(t, f_{\mathbf{w}})) = \partial \left( -t + \frac{1}{n(1-\epsilon)} \sum_{i=1}^n (t - \hat{\nu}_i)^+ \right) = -1 + \frac{1}{n(1-\epsilon)} \sum_{i=1}^{n(1-\epsilon)} \begin{cases} 1 & \text{if } t > \hat{\nu}_i \\ 0 & \text{if } t < \hat{\nu}_i \\ [0, 1] & \text{if } t = \hat{\nu}_i \end{cases} \quad (54)$$

We observe when setting  $t = \hat{\nu}_{n(1-\epsilon)}$ , it follows that  $0 \in \partial(-g(t, f_{\mathbf{w}}))$ . This is equivalent to the  $(1-\epsilon)$ -quantile of the Risk.  $\blacksquare$

### B.2 Proof of Lemma 4

**Proof.** By our choice of  $t^{(k+1)}$ , it follows:

$$\begin{aligned} \nabla_f g(t^{(k+1)}, f_{\mathbf{w}}^{(k)}) &= \nabla_f \left( t^{(k+1)} - \frac{1}{n(1-\epsilon)} \sum_{i=1}^n \left( t^{(k+1)} - \ell(\mathbf{x}_i; f_{\mathbf{w}}^{(k)}, y_i) \right)^+ \right) \\ &= -\frac{1}{n(1-\epsilon)} \sum_{i=1}^{n(1-\epsilon)} \nabla_f \left( t^{(k+1)} - \ell(\mathbf{x}_i; f_{\mathbf{w}}^{(k)}, y_i) \right)^+ = \frac{1}{n(1-\epsilon)} \sum_{i=1}^n \nabla_f \ell(\mathbf{x}_i; f_{\mathbf{w}}^{(k)}, y_i) \begin{cases} 1 & \text{if } t > \hat{\nu}_i \\ 0 & \text{if } t < \hat{\nu}_i \\ [0, 1] & \text{if } t = \hat{\nu}_i \end{cases} \end{aligned} \quad (55)$$

Now we note  $\nu_{n(1-\epsilon)} \leq t^{(k+1)} \leq \nu_{n(1-\epsilon)+1}$ . Then, plugging this into Equation (55), we have

$$\nabla_f g(t^{(k+1)}, f_{\mathbf{w}}^{(k)}) = \frac{1}{n(1-\epsilon)} \sum_{i=1}^{n(1-\epsilon)} \nabla_f \ell(\mathbf{x}_i; f_{\mathbf{w}}^{(k)}, y_i) \quad (56)$$

This concludes the proof.  $\blacksquare$

---

<sup>1</sup>In Progress.

### B.3 Projection onto a Norm Ball

In this section we show normalizing on to a norm-ball in the RKHS can be implemented efficiently.

**Lemma 14.** *Let  $\mathcal{K} \triangleq \{f : \|f\|_{\mathcal{H}} \leq R\}$ . Then, for a  $\hat{f} \notin \mathcal{K}$ , it follows*

$$\text{Proj}_{\mathcal{K}} \hat{f} = \left( \frac{R}{\|\hat{f}\|} \right) \hat{f} \quad (57)$$

**Proof.** We will formulate the dual problem and then find the corresponding  $f_{\mathbf{w}}$  that solves the dual.

$$\text{Proj}_{\mathcal{K}} \hat{f} = \arg \min_{f \in \mathcal{K}} \|f - \hat{f}\|_{\mathcal{H}}^2 = \arg \min_{f \in \mathcal{K}} \|f\|_{\mathcal{H}}^2 + \|\hat{f}\|_{\mathcal{H}}^2 - 2\langle f, \hat{f} \rangle_{\mathcal{H}} \quad (58)$$

$$= \arg \min_{f \in \mathcal{K}} \|f\|_{\mathcal{H}}^2 - 2\langle f, \hat{f} \rangle_{\mathcal{H}} \quad (59)$$

From here we can solve the dual problem. The Lagrangian is given by,

$$\mathcal{L}(f, u) \triangleq \|f\|_{\mathcal{H}}^2 - 2\langle f, \hat{f} \rangle + u \left( \|f\|_{\mathcal{H}}^2 - R^2 \right) \quad (60)$$

Then, we have dual problem as  $\theta(u) = \min_{f \in \mathcal{H}} \mathcal{L}(f, u)$ . Taking the derivative of the Lagrangian and setting it to zero, we obtain  $\arg \min_{f \in \mathcal{H}} \mathcal{L}(f, u) = (1 + u)^{-1} \hat{f}$ . With some more work, we obtain  $\arg \max_{u > 0} \theta(u) = R^{-1} \|\hat{f}\| - 1$ . We then have  $f$  at  $u^*$  as  $f = R \|\hat{f}\|_{\mathcal{H}}^{-1} \hat{f}$ . Since  $\|\hat{f}\| > R$  as  $\hat{f} \notin \mathcal{K}$  by assumption, our proof is complete. ■

## C Kernelized Binary Classification

In this section, we will prove error bounds for Subquantile Minimization in the Kernelized Binary Classification Problem.

### C.1 Subquantile Lipschitzness

**Lemma 15.** *( $L$ -Lipschitz of  $g(t, f)$  w.r.t  $f$ ). Let  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ , represent the data vectors. It then follows:*

$$|g(t, f) - g(t, \hat{f})| \leq L \|f - \hat{f}\|_{\mathcal{H}} \quad (61)$$

where  $L = \frac{1}{n(1-\epsilon)} \sum_{i \in X} \|\phi(\mathbf{x}_i)\|_{\mathcal{H}}$

**Proof.** We use the  $\mathcal{H}$  norm of the gradient to bound  $L$  from above. Let  $S$  be denoted as the subquantile set. Define the sigmoid function as  $\sigma(x) = \frac{1}{1+e^{-x}}$ .

$$\|\nabla_f g(t, f)\|_{\mathcal{H}} = \left\| \frac{1}{n(1-\epsilon)} \sum_{i=1}^n \mathbb{I} \left\{ t \geq (1 - y_i) \log(f^{(t)}(\mathbf{x}_i)) \right\} \left( y_i - \sigma(f^{(t)}(\mathbf{x}_i)) \right) \cdot \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \quad (62)$$

$$\stackrel{(i)}{\leq} \frac{1}{n(1-\epsilon)} \left\| \sum_{i \in S^{(t)}} \left( y_i - \sigma(f^{(t)}(\mathbf{x}_i)) \right) \cdot \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \quad (63)$$

$$\stackrel{(ii)}{\leq} \frac{1}{n(1-\epsilon)} \max_{i \in [n]} |y_i - \sigma(f^{(t)}(\mathbf{x}_i))| \left\| \sum_{i \in S^{(t)}} \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \quad (64)$$

$$\stackrel{(iii)}{\leq} \frac{1}{n(1-\epsilon)} \sum_{i \in X} \|\phi(\mathbf{x}_i)\|_{\mathcal{H}} \quad (65)$$

(i) follows from the triangle inequality. (ii) follows from the Cauchy-Schwarz inequality. (iii) follows from the fact that  $y_i \in \{0, 1\}$  and  $\text{range}(\sigma) \in [0, 1]$ . This completes the proof. ■

## C.2 Proof of Theorem 5

From Algorithm 1, we have for kernelized binary classification with linear kernel,

$$f^{(t+1)} = \text{Proj}_{\mathcal{K}} \left[ f^{(t)} - \frac{\eta}{n(1-\epsilon)} \sum_{i \in S^{(t)}} \left( \sigma(f^{(t)}(\mathbf{x}_i)) - y_i \right) \cdot \phi(\mathbf{x}_i) \right] \quad (66)$$

From which it follows,

$$\|f^{(t+1)} - f^*\|_{\mathcal{H}}^2 = \left\| \text{Proj}_{\mathcal{K}} \left[ f^{(t)} - \frac{\eta}{n(1-\epsilon)} \nabla g(f^{(t)}, t^*) \right] - f^* \right\|_{\mathcal{H}}^2 \quad (67)$$

$$\stackrel{(i)}{\leq} \left\| f^{(t)} - \frac{\eta}{n(1-\epsilon)} \nabla g(f^{(t)}, t^*) - f^* \right\|_{\mathcal{H}}^2 \quad (68)$$

$$= \|f^{(t)} - f^*\|_{\mathcal{H}}^2 - \frac{2\eta}{n(1-\epsilon)} \left\langle \nabla_f g(f^{(t)}, t^*), f^{(t)} - f^* \right\rangle_{\mathcal{H}} + \frac{\eta^2}{n^2(1-\epsilon)^2} \|\nabla_f g(f^{(t)}, t^*)\|_{\mathcal{H}}^2 \quad (69)$$

where (i) follows from the contraction property of the projection operator onto norm ball  $\mathcal{K}$ . We will expand the second term.

$$\frac{2\eta}{n(1-\epsilon)} \left\langle \nabla_f g(f^{(t)}, t^*), f^{(t)} - f^* \right\rangle_{\mathcal{H}} \stackrel{(66)}{=} \left\langle f^{(t)} - f^*, \frac{2\eta}{n(1-\epsilon)} \sum_{i \in S^{(t)}} \left( \sigma(f^{(t)}(\mathbf{x}_i)) - y_i \right) \cdot \phi(\mathbf{x}_i) \right\rangle_{\mathcal{H}} \quad (70)$$

$$= \left\langle f^{(t)} - f^*, \frac{2\eta}{n(1-\epsilon)} \sum_{i \in S^{(t)}} \left( \sigma(f^{(t)}(\mathbf{x}_i)) - \sigma(f^*(\mathbf{x}_i)) \right) \cdot \phi(\mathbf{x}_i) \right\rangle_{\mathcal{H}} \\ + \left\langle f^{(t)} - f^*, \frac{2\eta}{n(1-\epsilon)} \sum_{i \in S^{(t)}} (\sigma(f^*(\mathbf{x}_i)) - y_i) \cdot \phi(\mathbf{x}_i) \right\rangle_{\mathcal{H}} \quad (71)$$

We first upper bound upper bound the second term in Equation (71). From the Cauchy-Schwarz Inequality and noting  $y_i \in \{0, 1\}$  and  $\text{range}(\sigma) \in (0, 1)$ , we have the following,

$$\left\langle f^* - f^{(t)}, \frac{2\eta}{n(1-\epsilon)} \sum_{i \in S^{(t)}} (\sigma(f^*(\mathbf{x}_i)) - y_i) \phi(\mathbf{x}_i) \right\rangle_{\mathcal{H}} \quad (72)$$

$$\leq \frac{2\eta}{n(1-\epsilon)} \|f^* - f^{(t)}\|_{\mathcal{H}} \left\| \sum_{i \in S^{(t)}} \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \max_{i \in S^{(t)}} |\sigma(f^*(\mathbf{x}_i)) - y_i| \quad (73)$$

$$\stackrel{(ii)}{\leq} \frac{\eta^2}{n^{3/2}(1-\epsilon)^{3/2}} \|f^* - f^{(t)}\|_{\mathcal{H}}^2 \left\| \sum_{i \in S^{(t)}} \phi(\mathbf{x}_i) \right\|_{\mathcal{H}}^2 + \frac{2}{\sqrt{n(1-\epsilon)}} \sum_{i \in P} (\sigma(f^*(\mathbf{x}_i)) - y_i)^2 \quad (74)$$

where (ii) follows from Young's Inequality [You12] and noting for a vector  $\mathbf{x} \in \mathbb{R}^d$  it holds  $\|\mathbf{x}\|_{\infty} \leq \|\mathbf{x}\|_2$ . Let us now consider the function  $h : \mathcal{H} \rightarrow \mathbb{R}$  defined as  $h(f) \triangleq \sum_{i \in S \cap P} \log(1 + \exp(f^{(t)}(\mathbf{x}_i)))$ . We can then calculate the gradients by hand,  $\nabla h(f) = \sum_{i \in S \cap P} \sigma(f^{(t)}(\mathbf{x}_i)) \cdot \phi(\mathbf{x}_i)$  and  $\nabla^2 h(f) = \sum_{i \in S \cap P} \sigma(f^{(t)}(\mathbf{x}_i))(1 - \sigma(f^{(t)}(\mathbf{x}_i))) \cdot \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i)$ . From the Taylor Series expansion, we have for any  $f, \hat{f} \in \mathbb{R}^d$ , there exists  $\tilde{f}$  such that

$$\left\langle f - \hat{f}, \nabla h(f) - \nabla h(\hat{f}) \right\rangle_{\mathcal{H}} = \left\langle \nabla^2 h(\tilde{f}), (f - \hat{f}) \otimes (f - \hat{f}) \right\rangle_{\text{HS}} \quad (75)$$

Then, from the strong convexity of  $h$ , there exists a constant  $C$  such that the following inequality holds,

$$\left\langle f^{(t)} - f^*, \frac{2\eta}{n(1-\epsilon)} \sum_{i \in S^{(t)} \cap P} \left( \sigma(f^{(t)}(\mathbf{x}_i)) - \sigma(f^*(\mathbf{x}_i)) \right) \cdot \phi(\mathbf{x}_i) \right\rangle_{\mathcal{H}} \\ \gtrsim \frac{2\eta}{n(1-\epsilon)} \lambda_{\min} \left( \sum_{i \in S^{(t)}} \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) \right) \|f^{(t)} - f^*\|_{\mathcal{H}}^2 \stackrel{(iii)}{\geq} \frac{2\eta}{n(1-\epsilon)} \lambda_{\min} \left( \sum_{i \in S^{(t)} \cap P} \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) \right) \|f^{(t)} - f^*\|_{\mathcal{H}}^2 \quad (76)$$

where (iii) follows from Weyl's inequality [Wey12]. It is possible to show  $C = \Omega(\exp(-R \text{Tr}(\mathbf{\Sigma}) \log n))$  where  $R = \max_{t \in [T]} \|f^{(t)}\|_2$  from Proposition 9. We will now bound the final term in Equation (67).

$$\|\nabla_f g(f^{(t)}, t^*)\|_{\mathcal{H}}^2 = \left\| \frac{\eta}{n(1-\epsilon)} \sum_{i \in S^{(t)}} (\sigma(f^{(t)}(\mathbf{x}_i)) - y_i) \cdot \phi(\mathbf{x}_i) \right\|_{\mathcal{H}}^2 \quad (77)$$

$$\leq \frac{\eta^2}{n^2(1-\epsilon)^2} \max_{i \in S^{(t)}} |\sigma(f^{(t)}(\mathbf{x}_i)) - y_i|^2 \cdot \left\| \sum_{i \in S^{(t)}} \phi(\mathbf{x}_i) \right\|_{\mathcal{H}}^2 \quad (78)$$

$$\stackrel{(iv)}{\leq} \frac{\eta^2}{n^2(1-\epsilon)^2} \left\| \sum_{i \in S^{(t)}} \phi(\mathbf{x}_i) \right\|_{\mathcal{H}}^2 \cdot \sum_{i \in S^{(t)}} (\sigma(f^{(t)}(\mathbf{x}_i)) - y_i)^2 \quad (79)$$

$$\leq \frac{\eta^2}{n^2(1-\epsilon)^2} \left\| \sum_{i \in S^{(t)}} \phi(\mathbf{x}_i) \right\|_{\mathcal{H}}^2 \cdot \sum_{i \in P} (\sigma(f^{(t)}(\mathbf{x}_i)) - y_i)^2 \quad (80)$$

where (iv) follows from noting for any  $\mathbf{x} \in \mathbb{R}^d$  it holds  $\|\mathbf{x}\|_{\infty} \leq \|\mathbf{x}\|_{\mathcal{H}}$ . Now, combining Equations (67), (74), (76) and (77), we obtain

$$\begin{aligned} \|f^{(t+1)} - f^*\|_{\mathcal{H}}^2 &\leq \|f^{(t)} - f^*\|_{\mathcal{H}}^2 \left( 1 - \frac{2C\eta}{n(1-\epsilon)} \lambda_{\min} \left( \sum_{i \in S^{(t)} \cap P} \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) \right) + \frac{\eta}{2n(1-\epsilon)} \left\| \sum_{i \in S^{(t)}} \phi(\mathbf{x}_i) \right\|_{\mathcal{H}}^2 \right) \\ &\quad + \frac{\eta^2}{n^2(1-\epsilon)^2} \left\| \sum_{i \in S^{(t)}} \phi(\mathbf{x}_i) \right\|_{\mathcal{H}}^2 \cdot \sum_{i \in P} (\sigma(f^{(t)}(\mathbf{x}_i)) - y_i)^2 \end{aligned} \quad (81)$$

We will now expand out the final term in Equation (81) for  $(t+1)$ .

$$\frac{1}{n(1-\epsilon)} \sum_{i \in P} (\sigma(f^{(t+1)}(\mathbf{x}_i)) - y_i)^2 = \frac{1}{n(1-\epsilon)} \sum_{i \in P} (\sigma(f^{(t+1)}(\mathbf{x}_i)) - \sigma(f^*(\mathbf{x}_i)) + \sigma(f^*(\mathbf{x}_i)) - y_i)^2 \quad (82)$$

$$\leq \frac{2}{n(1-\epsilon)} \sum_{i \in P} (f^{(t+1)} - f^*)(\mathbf{x}_i)^2 + \frac{2}{n(1-\epsilon)} \sum_{i \in P} (\sigma(f^*(\mathbf{x}_i)) - y_i)^2 \quad (83)$$

$$= \frac{2}{n(1-\epsilon)} \left\langle \sum_{i \in P} \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i), (f^{(t+1)} - f^*) \otimes (f^{(t+1)} - f^*) \right\rangle_{\text{HS}} \quad (84)$$

$$\begin{aligned} &\quad + \frac{2}{n(1-\epsilon)} \sum_{i \in P} \left( \mathbf{Pr}_{(\mathbf{x}_i, y_i) \sim \mathcal{P}} \{y_i = +1 \mid \mathbf{x}_i\} - y_i \right)^2 \\ &\leq \frac{2}{n(1-\epsilon)} \|f^{(t+1)} - f^*\|_{\mathcal{H}}^2 \left\| \sum_{i \in P} \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} + \frac{2}{n(1-\epsilon)} \mathcal{E}_{\text{OPT}} \end{aligned} \quad (85)$$

Now, we can use Equation (81) to complete the bound.

$$\begin{aligned} &\frac{1}{n(1-\epsilon)} \sum_{i \in P} (\sigma(f^{(t+1)}(\mathbf{x}_i)) - y_i)^2 \\ &\leq \frac{2}{n(1-\epsilon)} \|f^{(t)} - f^*\|_{\mathcal{H}}^2 \left\| \sum_{i \in P} \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \left( 1 - \frac{2C\eta}{n(1-\epsilon)} \lambda_{\min}(\mathbf{X}_{\text{TP}}^{\top} \mathbf{X}_{\text{TP}}) + \frac{\eta^2}{n(1-\epsilon)} \left\| \sum_{i \in S^{(t)}} \phi(\mathbf{x}_i) \right\|_{\mathcal{H}}^2 \right) \\ &\quad + \frac{2}{n(1-\epsilon)} \mathcal{E}_{\text{OPT}} \end{aligned} \quad (86)$$

Now we define for all  $t \in [T]$ ,

$$\Lambda^{(t)} \triangleq \frac{1}{n(1-\epsilon)} \left\| \sum_{i \in P} \phi(\mathbf{x}_i) \otimes \phi(\mathbf{x}_i) \right\|_{\mathcal{H}} \|f^* - f^{(t)}\|_{\mathcal{H}}^2 + \frac{1}{n(1-\epsilon)} \sum_{i \in P} (\sigma(f^{(t)}(\mathbf{x}_i)) - y_i)^2 \quad (87)$$

We then have from Equations (77) and (86),

$$\Lambda^{(t+1)} \leq \max \left\{ 3 \left( 1 - \frac{2C\eta}{n(1-\epsilon)} \lambda_{\min}(\mathbf{X}_{\text{TP}}^\top \mathbf{X}_{\text{TP}}) + \frac{\eta^2}{n^{3/2}(1-\epsilon)^{3/2}} \left\| \sum_{i \in S^{(t)}} \phi(\mathbf{x}_i) \right\|_2^2 \right), \frac{\eta^2}{n^2(1-\epsilon)^2} \left\| \sum_{i \in S^{(t)}} \phi(\mathbf{x}_i) \right\|_{\mathcal{H}}^2 \right\} \cdot \Lambda^{(t)} + \frac{4}{\sqrt{n(1-\epsilon)}} \mathcal{E}_{\text{OPT}} \quad (88)$$

Solving the quadratic equation, we observe for sufficiently large  $n$  such that  $\frac{4C^2}{n^2(1-\epsilon)^2} \lambda_{\min}^2(\mathbf{X}_{\text{TP}}^\top \mathbf{X}_{\text{TP}}) \geq \frac{10}{3n(1-\epsilon)} \left\| \sum_{i \in S^{(t)}} \phi(\mathbf{x}_i) \right\|_2^2$  as the LHS scales in  $O(\lambda_{\min}(\mathbf{\Sigma}))$  and the RHS scales in  $O(\text{Tr}(\mathbf{\Sigma})/\sqrt{n})$ , and choosing

$$\eta \leq \frac{2C\lambda_{\min}(\mathbf{X}_{\text{TP}}^\top \mathbf{X}_{\text{TP}})}{\left\| \sum_{i \in S^{(t)}} \phi(\mathbf{x}_i) \right\|_2^2} = O\left(\frac{\lambda_{\min}(\mathbf{\Sigma})}{\text{Tr}(\mathbf{\Sigma})}\right) \quad (89)$$

We then observe there is a linear decrease in  $\Lambda^{(t)}$  plus a constant,

$$\Lambda^{(t+1)} \leq \max \left\{ \frac{1}{2}, \frac{4C^2\lambda_{\min}^2(\mathbf{X}_{\text{TP}}^\top \mathbf{X}_{\text{TP}})}{n^2(1-\epsilon)^2 \left\| \sum_{i \in S^{(t)}} \phi(\mathbf{x}_i) \right\|_2^2} \right\} \cdot \Lambda^{(t)} + \frac{4}{n(1-\epsilon)} \mathcal{E}_{\text{OPT}} \quad (90)$$

Then, noting  $\sum_{i \in P} (\sigma(f^{(t)}(\mathbf{x}_i)) - y_i)^2 \leq n(1-\epsilon)$ . We have  $\|f^{(t+1)} - f^*\|_{\mathcal{H}} \leq \varepsilon + O(\frac{\mathcal{E}_{\text{OPT}}}{\sqrt{n(1-\epsilon)}})$  after  $T = O\left(\log\left(\left(\frac{\lambda_{\max}(\mathbf{\Sigma})\|f^*\|_{\mathcal{H}}}{\sqrt{n}}\right) \frac{1}{\varepsilon}\right)\right)$  iterations. Our proof is complete.  $\blacksquare$

## D Proofs for Kernelized Multi-Class Classification

### D.1 Proof of Theorem 6

<sup>2</sup> From Algorithm 1, we have for kernelized multi-class classification with linear kernel,

$$f_k^{(t+1)} = \text{Proj}_{\mathcal{K}} \left[ f_k^{(t)} - \frac{\eta}{n(1-\epsilon)} \sum_{i \in S^{(t)}} \left( \text{softmax}_k(f^{(t)}(\mathbf{x}_i)) - y_{ik} \right) \cdot \mathbf{x}_i^\top \right] \quad (91)$$

where we adopt the notation of  $y_{ik}$  represents the  $k$ th element of  $\mathbf{y}_i$  which is equal to the canonical basis vector of  $\mathbb{R}^{|\mathcal{Y}|}$  with the label. From which it follows,

$$\|f^{(t+1)} - f^*\|_{\mathcal{H}}^2 = \sum_{k=1}^{|\mathcal{Y}|} \left\| \text{Proj}_{\mathcal{K}} \left[ f_k^{(t)} - \eta \nabla g(f_k^{(t)}, t^*) \right] - f_k^* \right\|_{\mathcal{H}}^2 \quad (92)$$

$$\stackrel{(i)}{\leq} \sum_{k=1}^{|\mathcal{Y}|} \|f_k^{(t)} - \eta \nabla g(f_k^{(t)}, t^*) - f_k^*\|_{\mathcal{H}}^2 \quad (93)$$

$$= \sum_{k=1}^{|\mathcal{Y}|} \|f_k^{(t)} - f_k^*\|_{\mathcal{H}}^2 - 2\eta \left\langle \nabla f g(f_k^{(t)}, t^*), f_k^{(t)} - f^* \right\rangle_{\mathcal{H}} + \eta^2 \|\nabla f g(f_k^{(t)}, t^*)\|_{\mathcal{H}}^2 \quad (94)$$

where (i) follows from the contraction property of the projection operator onto  $\mathcal{K}$ .

$$\left\langle \nabla f g(f_k^{(t)}, t^*), f_k^{(t)} - f_k^* \right\rangle_{\mathcal{H}} = \left\langle f_k^{(t)} - f^*, \sum_{i \in S^{(t)}} \left( \text{softmax}_k(f^{(t)}(\mathbf{x}_i)) - y_i \right) \cdot \mathbf{x}_i^\top \right\rangle_{\mathcal{H}} \quad (95)$$

$$\begin{aligned} &= \left\langle f_k^{(t)} - f^*, \sum_{i \in S^{(t)}} \left( \text{softmax}_k(f^{(t)}(\mathbf{x}_i)) - \text{softmax}_k(f^*(\mathbf{x}_i)) \right) \cdot \mathbf{x}_i^\top \right\rangle_{\mathcal{H}} \\ &\quad + \left\langle f_k^{(t)} - f^*, \sum_{i \in S^{(t)}} \left( \text{softmax}_k(f^*(\mathbf{x}_i)) - y_i \right) \cdot \mathbf{x}_i^\top \right\rangle_{\mathcal{H}} \end{aligned} \quad (96)$$

---

<sup>2</sup>In Progress



We will first lower bound the first term of Equation (96). Let us note the first fact that

$$\text{softmax}_k(f^{(t)}(\mathbf{x}_i)) - \text{softmax}_k(f^*(\mathbf{x}_i)) = \frac{\exp(f_k^{(t)}(\mathbf{x}_i))}{\zeta + \exp(f_k^{(t)}(\mathbf{x}_i))} - \frac{\exp(f^*(\mathbf{x}_i))}{\zeta + \exp(f^*(\mathbf{x}_i))} \quad (97)$$

where  $\zeta \triangleq \sum_{i=1, i \neq k}^{|\mathcal{Y}|} \exp(f_i^{(t)}(\mathbf{x}_i))$ . Define the function  $h(f) \triangleq \sum_{i \in S^{(t)}} \ln(\exp(f(\mathbf{x}_i)) + \zeta)$ . We can then manually compute the derivatives,  $\nabla h(f) \triangleq \sum_{i \in S^{(t)}} \frac{\exp(f(\mathbf{x}_i))}{\zeta + \exp(f(\mathbf{x}_i))} \cdot \mathbf{x}_i^\top$ , and  $\nabla^2 h(f) \triangleq \sum_{i \in S^{(t)}} \frac{\zeta \exp(f(\mathbf{x}_i))}{(\zeta + \exp(f(\mathbf{x}_i)))^2} \cdot \mathbf{x}_i \mathbf{x}_i^\top$ . Considering that  $f \in \mathcal{K}$  and  $\|\mathbf{x}_i\|$  is bounded almost surely, we have for a constant  $C$ ,

$$\begin{aligned} & \left\langle f_k^{(t)} - f_k^*, \frac{2\eta}{n(1-\epsilon)} \sum_{i \in S^{(t)}} \left( \text{softmax}_k(f^{(t)}(\mathbf{x}_i)) - \text{softmax}_k(f^*(\mathbf{x}_i)) \right) \cdot \mathbf{x}_i^\top \right\rangle_{\mathcal{H}} \\ & \gtrsim \frac{2\eta}{n(1-\epsilon)} \lambda_{\min} \left( \sum_{i \in S^{(t)}} \mathbf{x}_i \mathbf{x}_i^\top \right) \|f_k^{(t)}(\mathbf{x}_i) - f_k^*(\mathbf{x}_i)\|_{\mathcal{H}}^2 \stackrel{(i)}{\gtrsim} \frac{2\eta}{n(1-\epsilon)} \lambda_{\min} \left( \sum_{i \in S^{(t)} \cap P} \mathbf{x}_i \mathbf{x}_i^\top \right) \|f_k^{(t)}(\mathbf{x}_i) - f_k^*(\mathbf{x}_i)\|_{\mathcal{H}}^2 \end{aligned} \quad (98)$$

where (i) follows from Weyl's Inequality [Wey12] and from noting the function  $r(x) = \frac{x}{x+a}$  for  $a \in \mathbb{R}_+$  is monotone increasing over  $x \in \mathbb{R}_+$ . We will now upper bound the second term of Equation (96).

$$\begin{aligned} \left\langle f_k^{(t)} - f_k^*, \sum_{i \in S^{(t)}} (\text{softmax}_k(f^*(\mathbf{x}_i)) - y_{ik}) \right\rangle & \leq \|f_k^{(t)} - f_k^*\|_{\mathcal{H}} \left\| \sum_{i \in S^{(t)}} \mathbf{x}_i \right\|_2 \max_{i \in S^{(t)}} |\text{softmax}_k(f^*(\mathbf{x}_i)) - y_{ik}| \\ & \leq \frac{1}{2} \|f_k^{(t)} - f_k^*\|_{\mathcal{H}}^2 \left\| \sum_{i \in S^{(t)}} \mathbf{x}_i \right\|_2^2 + \left( \frac{1}{2} \wedge \sum_{i \in S^{(t)}} (\text{softmax}_k(f^*(\mathbf{x}_i)) - y_{ik})^2 \right) \end{aligned} \quad (99)$$

Now we will upper bound the final term in Equation (94).

$$\|\nabla g_k(f^{(t)}, t^*)\|_{\mathcal{H}}^2 = \left\| \sum_{i \in S^{(t)}} \left( \text{softmax}_k(f^{(t)}(\mathbf{x}_i)) - y_{ik} \right) \cdot \mathbf{x}_i \right\|_{\mathcal{H}}^2 \quad (100)$$

$$= \left\| \sum_{i \in S^{(t)}} \left( \text{softmax}_k(f^{(t)}(\mathbf{x}_i)) - \text{softmax}_k(f^*(\mathbf{x}_i)) + \text{softmax}_k(f^*(\mathbf{x}_i)) - y_{ik} \right) \cdot \mathbf{x}_i \right\|_{\mathcal{H}}^2 \quad (101)$$

$$\leq 2 \left\| \sum_{i \in S^{(t)}} \left( \text{softmax}_k(f^{(t)}(\mathbf{x}_i)) - \text{softmax}_k(f^*(\mathbf{x}_i)) \right) \right\|_{\mathcal{H}}^2 + 2 \left\| \sum_{i \in S^{(t)}} (\text{softmax}_k(f^*(\mathbf{x}_i)) - y_{ik}) \cdot \mathbf{x}_i \right\|_{\mathcal{H}}^2 \quad (102)$$