

# Hebbian learning inspired estimation of the linear regression parameters from queries

Johannes Schmidt-Hieber\* and Wouter M. Koolen†

September 26, 2023

## Abstract

Local learning rules in biological neural networks (BNNs) are commonly referred to as Hebbian learning. [26] links a biologically motivated Hebbian learning rule to a specific zeroth-order optimization method. In this work, we study a variation of this Hebbian learning rule to recover the regression vector in the linear regression model. Zeroth-order optimization methods are known to converge with suboptimal rate for large parameter dimension compared to first-order methods like gradient descent, and are therefore thought to be in general inferior. By establishing upper and lower bounds, we show, however, that such methods achieve near-optimal rates if only queries of the linear regression loss are available. Moreover, we prove that this Hebbian learning rule can achieve considerably faster rates than any non-adaptive method that selects the queries independently of the data.

**Keywords:** Biological neural networks, derivative-free methods, linear regression model, min-max estimation, sequential estimation, zeroth-order optimization.

**MSC 2020:** Primary: 62L20; secondary: 62J05

## 1 Introduction

While considerable efforts have been dedicated to develop the theoretical foundations underlying artificial neural networks (ANNs), a theory for learning in biological neural networks (BNNs) remains largely unexplored.

ANNs were designed to mimic BNNs but there are distinct differences in terms of the network structure and the learning. ANNs are deterministic and pass real numbers through the network whereas BNNs are stochastic and biological neurons send so-called spike trains. The information in BNNs is decoded in the spike times, the moments when the neurons fire/discharge. The brain updates the weights locally based on the state of the neurons that it connects. There

---

\*University of Twente

†University of Twente and CWI

are various quantitative and qualitative updating rules and they are commonly referred to as Hebbian learning. Specifically, spike-time dependent plasticity rules are Hebbian learning rules based on local spike times. It has been widely acknowledged in previous work that Hebbian learning in BNNs cannot perform gradient descent [10, 20, 29]. The rationale behind is that updating one parameter in a gradient descent scheme requires to evaluate a partial derivative that depends on the other parameters. This cannot be implemented by a local learning rule.

The recent work [26] derives a specific derivative-free optimization method that captures key characteristics of the spike time dependent plasticity rule underlying Hebbian learning in BNNs. For mathematical tractability, we consider here a slight variation of this scheme.

We study this iterative scheme in a scenario where instead of the full data, one can only query the linear regression model, see Section 2.2 for a definition and discussion. Based on  $k$  rounds of querying, the aim of the method is to recover the  $d$ -dimensional regression vector. Estimation in this setting is non-trivial. In particular, the data are not informative enough to compute gradients and run gradient descent. A first contribution of this article is to derive a bound for the convergence rate of the biologically inspired gradient-free learning rule in the query model. It is moreover shown that up to a log-factor in the number of parameters  $d$ , the derived convergence rates matches the lower bound for sample sizes  $k \gtrsim d^2 \log(d)$ . The derived minimax lower bound is non-standard. The main obstacle is that sequential estimation procedures with queries depending on previous observations induces dependence among the data. As we do not constraint the queries, the induced dependence on the future data is hard to characterize and to control. Compared to the earlier lower bounds for adaptive sensing [2] and derivative-free stochastic optimization [27, 12], the main difficulty here is that the dependence on previous queries and the parameter appears not in the mean but in the variance of the data distribution.

Any method achieving the optimal rate in the linear model with queries needs to carefully exploit the information obtained from previous queries. Indeed, we show that naive methods that specify the queries independently of the data will converge with a slower rate in the dimension parameter  $d$ , see Corollary 1 for the precise statement. This should be contrasted with adaptive sensing where it is known that the adaptation to previously seen data can improve the rate by, at most, a log-factor [2].

The article is structured as follows. Section 2 provides more details on biological learning and the link with gradient-free optimization. After formally stating the query model, the convergence rates are given. A corresponding lower bound is stated in Section 3. Section 4 derives matching upper and lower bounds for the convergence rate in the case of non-adaptive queries. This enables us to quantify the gain in the convergence rate by integrating previously observed query values into the query strategy. Related literature is summarized in Section 5. Proofs are deferred to the Appendix.

**Notation:** For vectors the spectral norm coincides with the Euclidean norm, see Lemma 4. Therefore, we can denote both norms by  $\|\cdot\|$ . Matrix inequalities are taken with respect to the partial ordering of symmetric matrices (Loewner order).

## 2 Upper bounds

### 2.1 Hebbian learning inspired update rule

Working in the supervised learning framework, suppose we want to learn a  $d$ -dimensional parameter vector  $\boldsymbol{\theta}$  from training data  $(\mathbf{X}_1, Y_1), (\mathbf{X}_2, Y_2), \dots$ , with  $d$ -dimensional design vectors/inputs  $\mathbf{X}_1, \mathbf{X}_2, \dots$  and real-valued response variables  $Y_1, Y_2, \dots$ . For our approach, it is sufficient to assume that the number of parameters and the number of covariates are the same, that is,  $\boldsymbol{\theta}$  and  $\mathbf{X}_1, \mathbf{X}_2, \dots$  are all vectors of length  $d$ .

Whereas gradient descent based methods are ubiquitous in machine learning, BNNs rely on local updating rules and receive feedback through neurotransmitters such as dopamine. Anticipating to solve a task well reduces the amount of released neurotransmitter in the brain. In reward-based synaptic plasticity, the amount of released neurotransmitter is modelled as the difference between a realized loss based on the current task and the anticipated loss that predicts the current loss based on previously seen losses [14].

Zeroth-order optimization methods are derivative-free methods that use the evaluation of the loss but do not involve the gradient of the loss. Since weight updates in BNNs are based on evaluations of the loss, it seems natural to interpret the learning in BNNs as a specific zeroth-order method. [26] extracts from the spike-time dependent local updating of the weights in a BNN a global zeroth-order rule for learning the parameter vector  $\boldsymbol{\theta}$ . If  $\boldsymbol{\theta}_k$  denotes the  $k$ -th update,  $L(\boldsymbol{\theta}, \mathbf{X}_k, Y_k)$  is the loss for parameter  $\boldsymbol{\theta}$  on the  $k$ -th training sample  $(\mathbf{X}_k, Y_k)$ ,  $\bar{L}_k$  is the anticipated loss in the  $k$ -th round based on previously seen losses, and  $\mathbf{U}_k$  is a  $d$ -dimensional uniform random vector  $\text{Unif}([A-, A]^d)$ , the biologically inspired update formula is

$$\boldsymbol{\theta}_k = \boldsymbol{\theta}_{k-1} + \alpha_k (L(\boldsymbol{\theta}_{k-1} + \mathbf{U}_k, \mathbf{X}_k, Y_k) - \bar{L}_k) (e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}). \quad (1)$$

The expressions  $e^{\pm \mathbf{U}_k}$  have to be understood componentwise, that is, if  $\mathbf{U}_k = (U_{k1}, \dots, U_{kd})^\top$ , then  $e^{\mathbf{U}_k} = (e^{U_{k1}}, \dots, e^{U_{kd}})^\top$  and  $e^{-\mathbf{U}_k} = (e^{-U_{k1}}, \dots, e^{-U_{kd}})^\top$ . The parameter  $A$  in the uniform distribution is a constant of the biological network but treated here as a hyperparameter of the optimization method. There is little guidance from the neuroscience on how to choose the anticipated loss  $\bar{L}_k$ . In [26] the  $\bar{L}_k$  was taken as the loss from the previous round  $L(\boldsymbol{\theta}_{k-2} + \mathbf{U}_{k-1}, \mathbf{X}_{k-1}, Y_{k-1})$ . For mathematical tractability, we assign instead the value  $L(\boldsymbol{\theta}_{k-1} + \mathbf{U}'_k, \mathbf{X}_k, Y_k)$  to  $\bar{L}_k$ , for an independently sampled and uniformly distributed random vector  $\mathbf{U}'_k \sim \text{Unif}([-A, A]^d)$ . Thus, we will study in this work convergence of the zeroth-order method

$$\boldsymbol{\theta}_k = \boldsymbol{\theta}_{k-1} + \alpha_k (L(\boldsymbol{\theta}_{k-1} + \mathbf{U}_k, \mathbf{X}_k, Y_k) - L(\boldsymbol{\theta}_{k-1} + \mathbf{U}'_k, \mathbf{X}_k, Y_k)) (e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}), \quad k = 1, \dots \quad (2)$$

Theorem 1 in [26] shows that in expectation, this rule does approximately gradient descent. Working here with a slightly different anticipated loss does not change this result. Indeed, since  $\mathbf{U}_k$  is independent of all other randomness and  $\mathbb{E}[e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}] = 0$ , conditioning on everything except for  $\mathbf{U}_k$  gives

$$\mathbb{E} \left[ L(\boldsymbol{\theta}_{k-1} + \mathbf{U}'_k, \mathbf{X}_k, Y_k) (e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}) \right] = \mathbb{E} \left[ L(\boldsymbol{\theta}_{k-1} + \mathbf{U}'_k, \mathbf{X}_k, Y_k) \mathbb{E}[e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}] \right] = 0.$$

Thus, as in the original version, the anticipated loss vanishes in expectation and Theorem 1 in [26] holds without any changes. While the connection to gradient descent is appealing, the main obstacle for fast convergence of zeroth-order methods is the high variance.

## 2.2 The query model

A natural first step of a statistical analysis of biologically inspired learning is to study the properties of the zeroth-order update rule (2) for standard statistical models. The overall aim is to identify relevant models where this rule achieves the optimal convergence rates and/or outperforms other standard methods.

A fundamental problem in statistics and machine learning is the case where we want to learn a regression/feature vector  $\boldsymbol{\theta}^*$  such that for an independent draw  $(\mathbf{X}, Y)$  with the same distribution as the (training) data  $(\mathbf{X}_k, Y_k)$ , we have  $Y \approx \mathbf{X}^\top \boldsymbol{\theta}^*$ . Considering squared loss

$$L(\boldsymbol{\theta}, \mathbf{X}, Y) = (Y - \mathbf{X}^\top \boldsymbol{\theta})^2,$$

the update formula (2) can then be rewritten as

$$\boldsymbol{\theta}_k = \boldsymbol{\theta}_{k-1} + \alpha_k \left( (Y_k - \mathbf{X}_k^\top (\boldsymbol{\theta}_{k-1} + \mathbf{U}_k))^2 - (Y_k - \mathbf{X}_k^\top (\boldsymbol{\theta}_{k-1} + \mathbf{U}'_k))^2 \right) (e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}). \quad (3)$$

Interestingly, this formula does not require full knowledge of the covariate vector  $\mathbf{X}_k$  and can be computed from queries. Here a query is defined as follows. Based on earlier seen queried observations and possibly extra randomness, the statistician chooses in the  $k$ -th round a query vector  $\mathbf{v}_k$ . Instead of the full  $k$ -th training sample  $(\mathbf{X}_k, Y_k)$ , one can only observe

$$Z_k = Y_k - \mathbf{X}_k^\top \mathbf{v}_k.$$

To realize (3), one needs to query each  $(\mathbf{X}_k, Y_k)$  twice. This means that in the  $k$ -th round, we choose two query vectors  $\mathbf{v}_k, \mathbf{v}'_k$  based on previously seen query values and extra randomness. As observations, we get

$$Z_k = Y_k - \mathbf{X}_k^\top \mathbf{v}_k, \quad \text{and} \quad Z'_k = Y_k - \mathbf{X}_k^\top \mathbf{v}'_k. \quad (4)$$

To see that this observational scheme is indeed enough to compute the updates (3), we can argue by induction, assuming that the initialization  $\boldsymbol{\theta}_0$  has been chosen independently of the data. The induction assumption is then, that  $\boldsymbol{\theta}_k$  only depends on the queries  $(Z_1, Z'_1, \dots, Z_k, Z'_k)$  and exogenous randomness. Then,  $\mathbf{v}_{k+1} = \boldsymbol{\theta}_k + \mathbf{U}_k$  and  $\mathbf{v}'_{k+1} = \boldsymbol{\theta}_k + \mathbf{U}'_k$  are admissible query vectors with corresponding queries  $Z_{k+1} = Y_{k+1} - \mathbf{X}_{k+1}^\top (\boldsymbol{\theta}_k + \mathbf{U}_{k+1})$  and  $Z'_{k+1} = Y_{k+1} - \mathbf{X}_{k+1}^\top (\boldsymbol{\theta}_k + \mathbf{U}'_{k+1})$ . The update formula implies that  $\boldsymbol{\theta}_{k+1}$  only depends on the queries  $(Z_1, Z'_1, \dots, Z_{k+1}, Z'_{k+1})$  and exogenous randomness, completing the induction step.

Queries can be thought of as an attention mechanism that tell us where to look next to extract useful information about the data. A query model seems appropriate if the full input vector cannot be processed because, for instance, the data arrive too quickly. The information from the queries is insufficient to compute gradients. Hence, gradient descent methods cannot be run in this case.

There is a wide variety of previously considered query models in the statistical literature. Examples are querying large graphs to avoid storage problems [21] or learning causal relationships from path queries [7].

To finish the description of the model, we have to specify the distribution of the latent/unobserved variables  $(\mathbf{X}_1, Y_1), (\mathbf{X}_2, Y_2), \dots$ . We assume that those are i.i.d. and generated from the linear regression model with random design, that is,  $\mathbf{X}_1, \mathbf{X}_2, \dots$  are drawn i.i.d. from some unknown distribution  $P_{\mathbf{X}}$  and the response variables are

$$Y_i = \mathbf{X}_i^\top \boldsymbol{\theta}^* + \varepsilon_i, \quad \text{for } i = 1, \dots \quad (5)$$

for i.i.d. noise variables  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$  that are also independent of the covariates  $\mathbf{X}_1, \dots$ . We assume that  $\sigma > 0$  is known.

Let  $(\mathbf{X}, Y)$  be a new and independent sample with the same distribution as the training samples. Write

$$Q := \mathbb{E}[\mathbf{X}\mathbf{X}^\top]$$

for the (uncentered) covariance matrix of the design vectors. Consider a possibly randomized estimator  $\hat{\boldsymbol{\theta}}$ . By construction,  $\hat{\boldsymbol{\theta}}$  is independent of a test point  $(\mathbf{X}, Y)$ . Rewriting  $Y - \mathbf{X}^\top \hat{\boldsymbol{\theta}} = (Y - \mathbf{X}^\top \boldsymbol{\theta}^*) + \mathbf{X}^\top (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)$ , conditioning on  $\hat{\boldsymbol{\theta}}$  and using tower rule, the excess risk of  $\hat{\boldsymbol{\theta}}$  is

$$\begin{aligned} \mathbb{E} \left[ (Y - \mathbf{X}^\top \hat{\boldsymbol{\theta}})^2 - (Y - \mathbf{X}^\top \boldsymbol{\theta}^*)^2 \right] &= \mathbb{E} \left[ 2(Y - \mathbf{X}^\top \boldsymbol{\theta}^*) \mathbf{X}^\top (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) + (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)^\top \mathbf{X} \mathbf{X}^\top (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \right] \\ &= \mathbb{E} \left[ (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)^\top Q (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) \right] \\ &= \text{tr}(QS) \end{aligned} \quad (6)$$

with  $S := \mathbb{E}[(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)^\top]$ .

[9] relates gradient descent with dropout to a vector autoregressive (VAR) process with random coefficients [24]. A similar argument can be made here. Defining  $\mathbf{W}_k := \boldsymbol{\theta}_k - \boldsymbol{\theta}^*$ ,  $G_k := I - 2\alpha_k(e^{-\mathbf{U}_k} - e^{\mathbf{U}_k})(\mathbf{U}'_k - \mathbf{U}_k)^\top \mathbf{X}_k \mathbf{X}_k^\top$  and  $\boldsymbol{\xi}_k := 2\alpha_k \varepsilon_k \mathbf{X}_k^\top (\mathbf{U}'_k - \mathbf{U}_k)(e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}) + \alpha_k((\mathbf{X}_k^\top \mathbf{U}_k)^2 - (\mathbf{X}_k^\top \mathbf{U}'_k)^2)(e^{-\mathbf{U}_k} - e^{\mathbf{U}_k})$ , the update formula can be written in the form of a lag one VAR process

$$\mathbf{W}_k = G_k \mathbf{W}_{k-1} + \boldsymbol{\xi}_k, \quad (7)$$

with independent random coefficients  $G_k$  and independent noise/innovation variables  $\boldsymbol{\xi}_k$ . It can be checked that the noise is centered, that is,

$$\mathbb{E}[\boldsymbol{\xi}_k] = 0. \quad (8)$$

A proof of this fact and a derivation of (7) can be found in Appendix A. As we can tune the learning rate  $\alpha_k$  and consider the parameter  $A$  from the uniform distribution as hyperparameter, it is interesting to work out the scaling of the different terms in these parameters. For small  $A$ ,  $e^{-\mathbf{U}_k} - e^{\mathbf{U}_k} \approx -2\mathbf{U}_k$  which is of order  $A$ . This means that the term  $2\alpha_k(e^{-\mathbf{U}_k} - e^{\mathbf{U}_k})(\mathbf{U}'_k - \mathbf{U}_k)^\top \mathbf{X}_k \mathbf{X}_k^\top$  in the definition of  $G_k$  will be of order  $\alpha_k A^2$ . The same is true for the first term of  $\boldsymbol{\xi}_k$  while the second term of  $\boldsymbol{\xi}_k$  scales like  $O(\alpha_k A^3)$ . This suggests to define

$$\text{effective learning rate} := \alpha_k A^2. \quad (9)$$

This quantity will remain constant if  $A$  is decreased by a factor  $\gamma < 1$ , and, concurrently,  $\alpha_k$  is increased by a factor  $\gamma^{-2}$ . Thus by making  $\gamma$  and thus  $A$  small, the term of the order  $O(\alpha_k A^3)$  becomes negligible. While the considered optimization scheme is motivated by Hebbian learning with  $A$  a small constant, treating  $A$  as a hyperparameter also suggests to choose a small  $A$ .

The updates (3) are independent of  $Q = \mathbb{E}[\mathbf{X}\mathbf{X}^\top]$ . As common in the literature on gradient descent, the chosen learning rate will, however, depend on the smallest eigenvalue of  $Q$ .

Applying (6) to the estimator  $\boldsymbol{\theta}_k$ , we need to control

$$S_k := \mathbb{E}[(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)^\top] = \mathbb{E}[\mathbf{W}_k \mathbf{W}_k^\top].$$

We can relate this to the previous iterate via  $\mathbf{W}_k \mathbf{W}_k^\top = G_k \mathbf{W}_{k-1} \mathbf{W}_{k-1}^\top G_k + \boldsymbol{\xi}_k \boldsymbol{\xi}_k^\top + G_k \mathbf{W}_{k-1} \boldsymbol{\xi}_k^\top + \boldsymbol{\xi}_k \mathbf{W}_{k-1}^\top G_k$ . Based on this identity, we will now derive a recursive formula for  $S_k$ . Assume that  $\mathbf{U}, \mathbf{U}' \sim \text{Unif}[-A, A]^d$  are independent and also independent of  $\mathbf{X}$ . Set  $\mathbf{D} = e^{-\mathbf{U}} - e^{\mathbf{U}}$  and let  $\mathbf{x}$  be a fixed  $d$ -dimensional vector. The following two matrices characterize the interaction between the  $d$ -dimensional noise vectors  $(\mathbf{U}, \mathbf{U}')$ ,

$$\begin{aligned} V(\mathbf{x}) &:= \mathbb{E}[\mathbf{x}^\top (\mathbf{U}' - \mathbf{U}) \mathbf{D} \mathbf{D}^\top (\mathbf{U}' - \mathbf{U}) \mathbf{x}] \\ W &:= \mathbb{E}[\left( (\mathbf{X}^\top \mathbf{U})^2 - (\mathbf{X}^\top \mathbf{U}')^2 \right)^2 \mathbf{D} \mathbf{D}^\top]. \end{aligned} \quad (10)$$

**Lemma 1.** (i) the interaction terms  $G_k \mathbf{W}_{k-1} \boldsymbol{\xi}_k^\top$  and  $\boldsymbol{\xi}_k \mathbf{W}_{k-1}^\top G_k$  have mean zero,

(ii)  $\text{Cov}(\boldsymbol{\xi}_k) = \mathbb{E}[\boldsymbol{\xi}_k \boldsymbol{\xi}_k^\top] = 4\alpha_k^2 \sigma^2 \mathbb{E}[V(\mathbf{X})] + \alpha_k^2 W$ ,

(iii) if  $\mu := -\mathbb{E}[U(e^{-U} - e^U)]$  with  $U \sim \text{Unif}([-A, A])$ , then,

$$\begin{aligned} S_k &= (I - 2\alpha_k \mu Q) S_{k-1} (I - 2\alpha_k \mu Q) + 4\alpha_k^2 \{ \mathbb{E}[\mathbf{X}^\top S_{k-1} \mathbf{X} V(\mathbf{X})] - \mu^2 Q S_{k-1} Q \} \\ &\quad + 4\alpha_k^2 \sigma^2 \mathbb{E}[V(\mathbf{X})] + \alpha_k^2 W. \end{aligned}$$

and moreover,

$$S_k \leq \|S_{k-1}\| \left( I - 4\alpha_k \mu Q + 4\alpha_k^2 \mathbb{E}[\|\mathbf{X}\|^2 V(\mathbf{X})] \right) + 4\alpha_k^2 \sigma^2 \mathbb{E}[V(\mathbf{X})] + \alpha_k^2 W. \quad (11)$$

In summary,  $S_k$  is an affine function in  $S_{k-1}$ .

As discussed before, the last term  $\alpha_k^2 W$  can be made negligible by choosing a sufficiently small parameter  $A$  in the uniform distribution. The term  $\mu$  is  $O(A^2)$  and, under moment assumptions on  $\mathbf{X}$ , the respective order of the terms  $\|\mathbb{E}[V(\mathbf{X})]\|$  and  $\|\mathbb{E}[\|\mathbf{X}\|^2 V(\mathbf{X})]\|$  is  $O(A^4 d)$  and  $O(A^4 d^2)$ . To decrease the norm of the error matrix  $S_k$  in (11), the factor  $I - 4\alpha_k \mu Q + 4\alpha_k^2 \mathbb{E}[\|\mathbf{X}\|^2 V(\mathbf{X})]$  has to remain below the identity matrix. This implies that  $\alpha_k A^2 Q \gtrsim \alpha_k^2 A^4 d^2 I$  and motivates to choose the effective learning rate  $\alpha_k A^2$  in (9) to be bounded by  $\lesssim \lambda_{\min}(Q)/d^2$ . Compared to standard choices such as  $1/k$ , the learning rate in the beginning is thus small and learning makes little progress.

If  $k \gtrsim d^2$ , effective learning rate  $1/k$  is possible by choosing  $\alpha_k A^2 \asymp 1/(k \vee \lambda_{\min}(Q) d^2)$ . The rate is then determined by the contributions  $4\alpha_k^2 \sigma^2 \|\mathbb{E}[V(\mathbf{X})]\| = O(\sigma^2 d/k^2)$ . The  $1/k^2$  becomes a  $1/k$  as the contributions from the previous  $O(k)$  iterates add up. Up to a  $\log^2(d)$ -factor, this explains the convergence rate stated in (13) below.

To simplify the exposition, we work now with specific choices for the parameter  $A$  in the uniform distribution and the learning rate  $\alpha_k$ . The expressions also depend on the variance of the noise  $\sigma^2$  in the linear regression model, which is assumed to be known. To obtain a convenient expressions for the upper bound, we moreover introduce

$$M_k := \max_{i=1, \dots, d} 1 \vee E[X_i^k],$$

where  $X_i$  denotes the  $i$ -th component of  $\mathbf{X}$ .

**Theorem 1.** *Assume  $d \geq 9$ . Let  $A = \sigma/\sqrt{d}$ , and choose the learning rate*

$$\alpha_k = \frac{11B \log(d)}{A^2 \lambda_{\min}(Q)(Bk + d^2 \log(d))} \quad \text{with } B := 1 \wedge \frac{\lambda_{\min}(Q)^2}{2904M_4}. \quad (12)$$

*Then there exists a constant  $C = C(M_4, \lambda_{\min}(Q))$ , such that for all  $k > 2d^2 \log(d)/B$ , we have*

$$\|\mathbb{E}[(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)^\top]\| \leq \left(\frac{2d^2 \log(d)}{Bk}\right)^{\log(d)} \|S_0\| + C \frac{\sigma^2 d \log^2(d)}{k}. \quad (13)$$

*If moreover  $\text{tr}(Q) \leq d$ , the excess risk (6) is bounded by*

$$\mathbb{E}[(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)^\top Q (\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)] \leq d \left(\frac{2d^2 \log(d)}{Bk}\right)^{\log(d)} \|S_0\| + C \frac{\sigma^2 d^2 \log^2(d)}{k}.$$

Assuming  $\text{tr}(Q) \leq d$  is natural and includes in particular the case that  $Q = \mathbb{E}[\mathbf{X}\mathbf{X}^\top]$  is the identity matrix.

Initializing  $\boldsymbol{\theta}_0 = \mathbf{0}$  with the zero vector, we have  $S_0 = \boldsymbol{\theta}^*(\boldsymbol{\theta}^*)^\top$  and with Lemma 4,  $\|S_0\| = \|\boldsymbol{\theta}^*\|^2$ .

The rate consists of two terms. For any  $\kappa > 0$ ,  $d > e^\kappa$ , and  $k \geq 2e^{\gamma+3\kappa}d^2 \log(d)/B$ , we have

$$\left(\frac{2d^2 \log(d)}{Bk}\right)^{\log(d)} \lesssim \frac{1}{d^\gamma k^\kappa}. \quad (14)$$

A proof of this inequality is given in Appendix A.4. The inequality states that for sufficiently large  $k \gtrsim d^2 \log(d)$ , the first term is negligible and the rate is  $\sigma^2 d^2 \log^2(d)/k$ . This rate is slower than the usual rate  $\sigma^2 d/k$  if the full data are observed in the linear regression model. The  $\log^2(d)$  factor seems to be an artifact of the proof and we will later derive nearly-matching lower bounds.

Why do we loose a factor  $d$  due to querying? It is tempting to assume that observed random variables from the query model and the linear regression model are equally informative. Then one can link the additional factor  $d$  for the query model to the loss of information in the data: While in the linear regression model we observe in every round a  $d$ -dimensional covariate vector  $\mathbf{X}_k$  together with a real-valued response  $Y_k$ , the query model observes in each round two real-valued queries. The total number of observed random variables in the query model is thus decreased by an order  $O(d)$ . Differently speaking the query model needs  $O(d)$  more iterations to receive the same number of observed variables. The concepts in [28] might be suitable to formalize such an argument.

The result is similar in spirit to the bounds obtained in [8] for another biologically motivated learning rule called (weight-perturbed) forward gradient descent [6, 25]. For a sequence of i.i.d.

random vectors  $\zeta_1, \zeta_2, \dots \mathcal{N}(0, I_d)$ , forward gradient descent with learning rate  $\alpha'_k$  is given by the update rule

$$\boldsymbol{\theta}_k = \boldsymbol{\theta}_{k-1} - \alpha'_k (\nabla L(\boldsymbol{\theta}_{k-1}))^\top \zeta_k \zeta_k, \quad k = 1, 2, \dots$$

The similarity is best explained for squared loss  $L$ . In this case,  $\nabla L(\boldsymbol{\theta}) = -2(Y_k - \mathbf{X}_k^\top \boldsymbol{\theta}) \mathbf{X}_k$ . If instead of an independent draw, we set  $\mathbf{U}'_k := -\mathbf{U}_k$ , then the optimization scheme (3) becomes

$$\begin{aligned} \boldsymbol{\theta}_k &= \boldsymbol{\theta}_{k-1} + \alpha_k \left( (Y_k - \mathbf{X}_k^\top (\boldsymbol{\theta}_{k-1} + \mathbf{U}_k))^2 - (Y_k - \mathbf{X}_k^\top (\boldsymbol{\theta}_{k-1} - \mathbf{U}_k))^2 \right) (e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}) \\ &= \boldsymbol{\theta}_{k-1} + 2\alpha_k \nabla L(\boldsymbol{\theta}_{k-1})^\top \mathbf{U}_k (e^{\mathbf{U}_k} - e^{-\mathbf{U}_k}). \end{aligned}$$

For small  $A$ ,  $e^{\mathbf{U}_k} - e^{-\mathbf{U}_k} \approx 2\mathbf{U}_k$ , which means that  $\boldsymbol{\theta}_k \approx \boldsymbol{\theta}_{k-1} + 4\alpha_k \nabla L(\boldsymbol{\theta}_{k-1})^\top \mathbf{U}_k \mathbf{U}_k$ . Whereas forward gradient descent has been proposed for normally distributed  $\zeta_k$ , the gradient-free optimization scheme here is closely related to the case where these vectors are sampled from a uniform distribution. The choice  $\mathbf{U}'_k := -\mathbf{U}_k$  leads to less remainder terms in the analysis and less additional noise, but we find this choice less appealing to model Hebbian learning in the brain.

### 3 Lower bound for adaptive queries

We now derive nearly-matching lower bounds. To make the lower bound rigorous, we need to formalize the query model. The observed data consists of  $2k$  query vectors  $\mathbf{v}_1, \mathbf{v}'_1, \dots, \mathbf{v}_k, \mathbf{v}'_k$  and a random vector  $(Z_1, Z'_1, \dots, Z_k, Z'_k)$  of length  $2k$ .

The unobserved/latent i.i.d. pairs  $(\mathbf{X}_\ell, Y_\ell)$ ,  $\ell = 1, \dots, k$  are generated from the linear regression model (5). For the lower bounds we will moreover assume that the design distribution  $P_{\mathbf{X}}$  is  $\mathcal{N}(0, I_d)$  with  $I_d$  the  $d \times d$  identity matrix. This implies that  $Q = E[\mathbf{X}\mathbf{X}^\top] = I_d$ .

Let  $(\mathcal{G}_\ell)_{\ell \geq 0}$  be a filtration generated by exogenous randomness. We assume that for any  $\ell = 1, 2, \dots$ , the query vectors  $\mathbf{v}_\ell, \mathbf{v}'_\ell$  are measurable with respect to the  $\sigma$ -algebra

$$\mathcal{F}_{\ell-1} := \sigma(Z_1, Z'_1, \dots, Z_{\ell-1}, Z'_{\ell-1}, \mathbf{v}_1, \mathbf{v}'_1, \dots, \mathbf{v}_{\ell-1}, \mathbf{v}'_{\ell-1}) \times \mathcal{G}_{\ell-1}.$$

This means that the query vectors might depend on past queries, past query vectors and exogenous randomness. Given  $\mathbf{v}_\ell, \mathbf{v}'_\ell$ , we observe in the  $\ell$ -th step the queries  $(Z_\ell, Z'_\ell)$  with

$$Z_\ell = Y_\ell - \mathbf{X}_\ell^\top \mathbf{v}_\ell, \quad Z'_\ell = Y_\ell - \mathbf{X}_\ell^\top \mathbf{v}'_\ell. \quad (15)$$

Finally, an estimator  $\hat{\boldsymbol{\theta}}_k$  is any measurable function in  $\mathcal{F}_k$ . This means that an estimator is allowed to depend on all observed queries, all observed query vectors and the exogenous randomness.

To allow for two queries instead of one makes the lower bounds considerable more involved as  $Z_\ell$  and  $Z'_\ell$  are dependent by virtue of sharing  $\mathbf{X}_\ell, Y_\ell$ . One query is, however, not enough to realize the zeroth-order scheme (3).

Before stating the lower bounds, we derive an equivalent representation of this model. Two statistical models  $(P_\boldsymbol{\theta} : \boldsymbol{\theta} \in \Theta)$  and  $(Q_\boldsymbol{\theta} : \boldsymbol{\theta} \in \Theta)$  with the same parameter space  $\Theta$  are equivalent if there exist Markov kernels  $M, M'$  that are independent of unknown parameters such that  $Q_\boldsymbol{\theta} = MP_\boldsymbol{\theta}$  and  $P_\boldsymbol{\theta} = M'Q_\boldsymbol{\theta}$  for all  $\boldsymbol{\theta} \in \Theta$ .



**Lemma 2.** Consider the query model  $\{Z_1, Z'_1, \dots, Z_k, Z'_k, \mathbf{v}_1, \mathbf{v}'_1, \dots, \mathbf{v}_k, \mathbf{v}'_k\}$  with (15) replaced by

$$Z_\ell = Y_\ell - \mathbf{X}_\ell^\top \mathbf{v}_\ell, \quad Z'_\ell = \mathbf{X}_\ell^\top \mathbf{v}'_\ell, \quad \ell = 1, 2, \dots \quad (16)$$

Both models are statistically equivalent.

For the lower bounds we will work in the transformed query model (16). If we choose the query vector  $\mathbf{v}'_\ell = 0$ , then  $Z'_\ell = 0$  and this is as informative as only querying the model once.

For the model with queries (16), we have

$$\begin{pmatrix} Z_\ell \\ Z'_\ell \end{pmatrix} \Big|_{\mathbf{v}_\ell, \mathbf{v}'_\ell} \sim \mathcal{N} \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma^2 + \|\boldsymbol{\theta} - \mathbf{v}_\ell\|^2 & \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \\ \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle & \|\mathbf{v}'_\ell\|^2 \end{pmatrix} \right). \quad (17)$$

using that, we also assumed  $\mathbf{X}_k \sim \mathcal{N}(0, I)$  and thus  $Q = I_d$ .

The choice of the query vectors up to round  $k$ , will be called the *query strategy* (up to round  $k$ ) and is denoted by  $\mathcal{V}_k$ . It is important to recall that the query strategy cannot depend on  $\boldsymbol{\theta}$  as it is unknown.

The data distribution of  $(Z_1, Z'_1, \dots, Z_k, Z'_k, \mathbf{v}_1, \mathbf{v}'_1, \dots, \mathbf{v}_k, \mathbf{v}'_k)$  will be denoted by  $P_{\boldsymbol{\theta}, \mathcal{V}_k}$  and depends on the underlying parameter  $\boldsymbol{\theta}$  and the query strategy  $\mathcal{V}_k$ . Since the choice of the query strategy  $\mathcal{V}_k$  is part of the estimation procedure, the right notion of the minimax estimation risk is

$$\inf_{\hat{\boldsymbol{\theta}}, \mathcal{V}_k} \sup_{\boldsymbol{\theta} \in \Theta} \mathbb{E}_{\boldsymbol{\theta}, \mathcal{V}_k} [\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|^2],$$

where the infimum is taken over all query strategies and all estimators. The parameter space  $\Theta$  will be chosen as the Euclidean ball with radius  $R$ , that is,

$$B_R(0) := \{\boldsymbol{\theta} : \|\boldsymbol{\theta}\| \leq R\}.$$

**Theorem 2.** If  $d \geq 3$  and  $k \geq d^2$ , then,

$$\inf_{\hat{\boldsymbol{\theta}}, \mathcal{V}_k} \sup_{\boldsymbol{\theta} \in B_R(0)} \mathbb{E}_{\boldsymbol{\theta}, \mathcal{V}_k} [\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|^2] \geq \frac{1}{162} \left(1 - \frac{1}{\sqrt{2}}\right) \left(R^2 \wedge \frac{d^2}{k} \sigma^2\right).$$

The statement assumes that  $k \geq d^2$ . A similar condition ( $k > 2d^2 \log(d)/B$ ) appears also in the corresponding upper bound in Theorem 1.

Interestingly, the rate does not depend on the radius  $R$  for large  $k$ . This indicates that the statistical main difficulty is to recover the direction of the true regression vector.

The same upper and lower bounds hold if instead of (4), we observe the squared queries  $(Z_1^2, (Z'_1)^2, \dots, Z_k^2, (Z'_k)^2)$ . Indeed the squares are sufficient to implement the biologically inspired updating rule (3). Since  $(Z_1^2, (Z'_1)^2, \dots, Z_k^2, (Z'_k)^2)$  is at most as informative as  $(Z_1, Z'_1, \dots, Z_k, Z'_k)$ , the lower bounds remain true.

## 4 Minimax risk for nonadaptive queries

What is the advantage to use previous information to select the next query vector? To answer this, we now consider query vectors  $\mathbf{v}_1, \dots, \mathbf{v}_k$  that are chosen before the data are revealed. Denote by  $\mathcal{M}_k$  the space of all such query strategies. For convenience, we will moreover assume throughout this section that the design distribution is standard multivariate normal, this means that

$$\mathbf{X}_1, \mathbf{X}_2, \dots \sim P_{\mathbf{X}} = \mathcal{N}(0, I_d), \quad \text{i.i.d.} \quad (18)$$

We show in this section that the minimax estimation risk with non-adaptive query strategies in  $\mathcal{M}_k$  and parameter space  $\Theta$  the Euclidean ball  $B_R(0)$  is

$$\inf_{\hat{\boldsymbol{\theta}}, \mathcal{V}_k \in \mathcal{M}_k} \sup_{\boldsymbol{\theta} \in B_R(0)} \mathbb{E}_{\boldsymbol{\theta}, \mathcal{V}_k} [\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|^2] \asymp R^2 \wedge \frac{d^2}{k} (R \vee \sigma)^2. \quad (19)$$

For the upper bound we construct the following estimator. If  $k \leq 2d^2((\sigma^2/R^2) \vee 1)$ , take  $\hat{\boldsymbol{\theta}} = 0$ . If  $k > 2d^2((\sigma^2/R^2) \vee 1)$ , we have  $d(k/(2d) + 1) = k/2 + d \leq k$  and can therefore partition the index set  $\{1, \dots, k\}$  into  $d$  blocks  $\mathcal{B}_1, \dots, \mathcal{B}_d$  such that each block has cardinality  $\geq k/(2d)$ . Set  $\mathbf{v}_j = (R \vee \sigma)\mathbf{e}_s$  if  $j \in \mathcal{B}_s$  with  $\mathbf{e}_s$  the  $s$ -th standard basis vector. Thanks to (18), the data are then given by

$$Z_j \sim \mathcal{N}(0, \sigma^2 + \|\boldsymbol{\theta} - (R \vee \sigma)\mathbf{e}_s\|^2) = \mathcal{N}(0, \sigma^2 + \|\boldsymbol{\theta}\|^2 - 2(R \vee \sigma)\boldsymbol{\theta}_s + (R \vee \sigma)^2), \quad \text{if } j \in \mathcal{B}_s$$

and the estimator for the  $s$ -th component of  $\boldsymbol{\theta}$  is in this case

$$\hat{\theta}_s := \frac{\sigma^2 + \|\boldsymbol{\theta}\|^2 + (R \vee \sigma)^2 - |\mathcal{B}_s|^{-1} \sum_{r \in \mathcal{B}_s} Z_r^2}{2(R \vee \sigma)}, \quad \text{for } s = 1, \dots, d, \quad (20)$$

with  $|\mathcal{B}_s|$  the cardinality of the set  $\mathcal{B}_s$ . Whenever  $r \in \mathcal{B}_s$ , we have  $\mathbb{E}[Z_r^2] = \sigma^2 + \|\boldsymbol{\theta}\|^2 - 2(R \vee \sigma)\boldsymbol{\theta}_s + (R \vee \sigma)^2$ , implying that

$$\hat{\boldsymbol{\theta}} = (\hat{\theta}_1, \dots, \hat{\theta}_d)$$

is an unbiased estimator for  $\boldsymbol{\theta}$ .

**Theorem 3.** *Assume (18). For the estimator  $\hat{\boldsymbol{\theta}}$  defined componentwise in (20), we have*

$$\sup_{\boldsymbol{\theta} \in B_R(0)} \mathbb{E}_{\boldsymbol{\theta}, \mathcal{V}_k} [\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|^2] \leq 25 \left( R^2 \wedge \frac{d^2}{k} (R \vee \sigma)^2 \right).$$

The estimator uses knowledge of  $R$  and  $\sigma^2$ . If these quantities are unknown, estimation seems unequal harder, in particular if  $R$  is small.

The query vectors  $\mathbf{v}_j$  can be thought of as test functions or features for  $\boldsymbol{\theta}$ . In the construction of the estimator, they have norm  $\|\mathbf{v}_j\| = (R \vee \sigma)$ . Interestingly, if  $\sigma > R$ , the norm exceeds  $R$  and the  $\mathbf{v}_j$  are themselves not in the parameter space  $\Theta = B_R(0)$ .

*Proof.* For  $k \leq 2d^2((\sigma^2/R^2) \vee 1)$ ,  $\hat{\boldsymbol{\theta}} = 0$  and the result follows since  $\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|^2 = \|\boldsymbol{\theta}\|^2 \leq R^2$ .

It remains to show that for  $k > 2d^2((\sigma^2/R^2) \vee 1)$ , the rate is bounded by  $25\frac{d^2}{k}(R \vee \sigma)^2$ . The bias-variance decomposition yields

$$\mathbb{E}_{\boldsymbol{\theta}, \nu_k} [\|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|^2] = \sum_{s=1}^d \mathbb{E}_{\boldsymbol{\theta}, \nu_k} [(\widehat{\boldsymbol{\theta}}_s - \boldsymbol{\theta}_s)^2] = \sum_{s=1}^d \text{Bias}_{\boldsymbol{\theta}, \nu_k}^2(\widehat{\boldsymbol{\theta}}_s) + \sum_{s=1}^d \text{Var}_{\boldsymbol{\theta}, \nu_k}(\widehat{\boldsymbol{\theta}}_s).$$

As we have already shown that  $\widehat{\boldsymbol{\theta}}$  is unbiased, it remains to bound the variances. For  $\xi \sim \mathcal{N}(0, a^2)$ , we have  $\text{Var}(\xi^2) = \mathbb{E}[\xi^4] - \mathbb{E}^2[\xi^2] = 3a^4 - a^2 = 2a^2$ . Using the definition of the estimator in (20), the independence of  $Z_1, \dots, Z_k$ , that  $\|\boldsymbol{\theta}\|^2 \leq R^2$ , that  $\sigma^2 + \|\boldsymbol{\theta}\|^2 - 2R\boldsymbol{\theta}_s + R^2 \leq 5(R \vee \sigma)^2$ , and that by construction of the blocks  $|\mathcal{B}_s| \geq k/(2d)$ , we obtain for any  $j \in \mathcal{B}_s$ ,

$$\text{Var}_{\boldsymbol{\theta}, \nu_k}(\widehat{\boldsymbol{\theta}}_s) = \frac{\text{Var}(Z_j^2)}{(2(R \vee \sigma))^2 |\mathcal{B}_s|} = \frac{2(\sigma^2 + \|\boldsymbol{\theta}\|^2 - 2R\boldsymbol{\theta}_s + R^2)^2}{4(R \vee \sigma)^2 |\mathcal{B}_s|} \leq \frac{50(R \vee \sigma)^4}{4(R \vee \sigma)^2 |\mathcal{B}_s|} \leq \frac{25d}{k}(R \vee \sigma)^2.$$

Summing over  $s = 1, \dots, d$  gives another factor  $d$  and thus the claim follows.  $\square$

We now state the corresponding lower bound.

**Theorem 4.** *Assume (18). If  $d \geq 6$ , then for any  $k = 1, 2, \dots$*

$$\inf_{\widehat{\boldsymbol{\theta}}, \nu_k \in \mathcal{M}_k} \sup_{\boldsymbol{\theta} \in B_{R(0)}} \mathbb{E}_{\boldsymbol{\theta}, \nu_k} [\|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|^2] \geq 2^{-18} \left( R^2 \wedge \frac{d^2}{k}(R \vee \sigma)^2 \right).$$

Together with Theorem 3, this shows (19). While in the adaptive setting, we had to impose the restriction  $k \gtrsim d^2 \log(d)$  for the upper bound and  $k \geq d^2$  for the lower bound, the derived rate in the non-adaptive setting holds for all sample sizes  $k$ .

Compared to the adaptive case, proving the lower bound in the non-adaptive case is considerably more involved. Reasons are that in this case, one additional regime occurs in the rate. To deal with this regime requires to show that whatever the choice of the query vectors is, one can find a regression vector  $\boldsymbol{\theta}^*$  in the parameter space that is far away. The fact that we allow for two queries per sample  $(\mathbf{X}_\ell, Y_\ell)$  makes that ‘far away’ has to be interpreted with respect to some tube that is generated by the pair of query vectors  $(\mathbf{v}_\ell, \mathbf{v}'_\ell)$ .

We believe that the constant  $2^{-18}$  in the lower bound can be improved significantly at the expense of a more technical proof.

The upper bound only needs one query per iteration and the lower bound is derived for two queries per iteration. This already proves that if we can only query once in every iteration, the minimax rate remains the same.

Based on the derived lower bound, we can now quantify the gap between adaptive and non-adaptive design. A natural setting is to allow that all parameters are of order one. This means that  $\|\boldsymbol{\theta}\|^2 = \sum_{j=1}^d \theta_j^2$  is of order  $O(d)$  and motivates to choose  $R = \sqrt{d}$ .

**Corollary 1.** *Assume (18). If  $\sigma = 1$ ,  $Q = I$ ,  $\boldsymbol{\theta}_0 = 0$ , and  $R = \sqrt{d}$ , then for all  $k \geq 2e^5 d^2 \log(d)/B$  (with  $B$  the constant in (12)) and  $d \geq 8$ , the upper for the adaptive design yields the convergence rate*

$$\frac{d^2 \log^2(d)}{k}$$

while the minimax rate for the non-adaptive design is

$$\frac{d^3}{k}.$$

The result implies that constructing queries based on previously seen data improves the rate by a factor  $d^{-1}$  (up to logarithms). The improvement will become even more pronounced if  $R$  increases. Indeed, the upper bound in the adaptive query setting will remain  $d^2 \log^2(d)/k$ , while the minimax rate for the non-adaptive query setting becomes  $d^2 R^2/k$ . We do not have a convincing heuristic argument explaining the gap in the rates. However, it is clear that the adaptive query setting can learn over time about the direction of the true  $\boldsymbol{\theta}^*$ , whereas in the non-adaptive query setting one has to spread out the query vectors equally over all possible directions. This is also clearly visible in the construction of the estimator in (20).

For the related problem of adaptive sensing, it has been found in [2] that adaptation improves the rate by at most log-factors. In this setting there are no queries and one can choose in the  $i$ -th iteration a design vector  $\mathbf{X}_i$  based on past observations and will then observe  $Y_i = \mathbf{X}_i^\top \boldsymbol{\theta}^* + \varepsilon_i$  with independent  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ . If  $s$  is the number of non-zero components of  $\boldsymbol{\theta}^*$ , it is shown that the risk  $\mathbb{E}[\|\hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*\|^2]$  of any estimator  $\hat{\boldsymbol{\theta}}$  based on  $k$  measurements is lower bounded by  $\gtrsim \sigma^2 s/k$ . On the contrary, in the non-adaptive setting with  $\mathbf{X}_i$  chosen i.i.d. and independent of previous data, the Dantzig selector  $\hat{\boldsymbol{\theta}}_k^D$  based on  $k \gtrsim s \log(d/s)$  measurements achieves estimation rate  $\mathbb{E}[\|\hat{\boldsymbol{\theta}}_k^D - \boldsymbol{\theta}^*\|^2] \lesssim \sigma^2 s \log(d)/k$ . Thus the gain in the convergence rate of an adaptive sampling strategy is here at most a factor  $\log(d)$  in the rate. The reason why adaptation hardly improves the rate in this setting is attributed in [2] to the difficulty to recover the support of the sparse regression vector  $\boldsymbol{\theta}^*$ .

## 5 Related literature

In zeroth-order stochastic convex optimization the task is to learn a minimizer of an unknown convex function  $f(\boldsymbol{\theta}) = \mathbb{E}_\Xi[f(\boldsymbol{\theta}, \Xi)]$ . The vanilla framework is to sequentially issue queries  $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots$  and receive noisy observations  $f(\boldsymbol{\theta}_k, \Xi_k)$  for i.i.d. unobserved  $\Xi_1, \Xi_2, \dots$ . The particular case of linear regression  $Y = \mathbf{X}^\top \boldsymbol{\theta}^* + \epsilon$  corresponds to  $\Xi = (\mathbf{X}, \epsilon)$  with  $d$ -dimensional covariate vectors  $\mathbf{X}$  drawn from an unknown distribution, noise  $\epsilon$ , and feedback  $f(\boldsymbol{\theta}, (\mathbf{X}, \epsilon)) = (\mathbf{X}^\top (\boldsymbol{\theta}^* - \boldsymbol{\theta}) + \epsilon)^2 = (Y - \mathbf{X}^\top \boldsymbol{\theta})^2$ . In the literature (we refer to [19, 5] for excellent surveys), rates for algorithms and lower bounds are organised based on three main distinctions:

- in the *optimization* literature the objective is the gap  $f(\hat{\boldsymbol{\theta}}_k) - \min_{\boldsymbol{\theta}} f(\boldsymbol{\theta})$  of a proposed evaluation point  $\hat{\boldsymbol{\theta}}_k$ , while in the *bandit* literature the objective is the regret of the queries issued  $\sum_{t=1}^k (f(\boldsymbol{\theta}_t) - \min_{\boldsymbol{\theta}} f(\boldsymbol{\theta}))$ . The latter bounds the former by online to batch conversion, but there are interesting separations in minimax rates [27].
- in *one-point feedback* the learner issues a point  $\boldsymbol{\theta}$  and observes  $f(\boldsymbol{\theta}, \Xi)$ , where the noise  $\Xi$  is i.i.d. between queries [13]. In *two-point feedback* the learner issues a pair of points  $\boldsymbol{\theta}, \boldsymbol{\theta}'$  and observes  $f(\boldsymbol{\theta}, \Xi)$  and  $f(\boldsymbol{\theta}', \Xi)$  with shared noise  $\Xi$ . Different rates can occur between one and two-point feedback.

- in addition to convexity of  $f$ , one may get better rates by assuming that  $f(\cdot)$  (or  $f(\cdot, \xi)$  for every  $\xi$ ) are Lipschitz, smooth, higher-order  $\beta$ -smooth [4], and/or strongly convex.

In [12], upper and lower bounds are derived for the gap  $f(\widehat{\boldsymbol{\theta}}_k) - \min_{\boldsymbol{\theta} \in \Theta} f(\boldsymbol{\theta})$ . Under the imposed conditions and ignoring the dependence on the number of parameters  $d$ , the gap decreases in  $k$  with the rate  $1/\sqrt{k}$ , which is slower than the  $1/k$  rate obtained here. The lower bounds are obtained for linear  $f(\boldsymbol{\theta}) = \mathbb{E}[\boldsymbol{\theta}^\top \mathbf{X}] = \boldsymbol{\theta}^\top \mathbb{E}[\mathbf{X}]$  with unknown distribution  $\mathbf{X}$ . In this case, the Hessian is zero and the minimizer will lie at the boundary of the parameter space  $\Theta$ . Therefore the rates strongly depend on the choice of  $\Theta$ .

In the case of linear regression, the function  $f(\boldsymbol{\theta}) = E[(Y - \mathbf{X}^\top \boldsymbol{\theta})^2]$  is a convex quadratic defined on the entire Euclidean space  $\mathbb{R}^d$ , and as such strongly convex and smooth (of infinite order) but not Lipschitz. For one-point feedback, [22] proves that, under suitable assumptions, the averaged iterations  $\bar{\boldsymbol{\theta}}_k = \frac{1}{k} \sum_{\ell=1}^k \boldsymbol{\theta}_\ell$  converge with rate  $\frac{d^2}{k} (\frac{k}{d})^{1/\beta}$  for any  $\beta > 0$ .

Keeping the dimension  $d$  fixed and letting the number of iterations tend to infinity, [16] derives a CLT for the average over all iterates (Ruppert-Polyak average).

To complete this literature overview, we briefly mention related approaches. [27] considers a zeroth-order method to learn a minimizer of the function  $F(\mathbf{w}) = \mathbf{w}^\top A \mathbf{w} + \mathbf{b}^\top \mathbf{w} + c$  for unknown  $d \times d$  matrix  $A$ ,  $d$ -dimensional vector  $\mathbf{b}$  and scalar  $c$ . In every iteration, one can query the function  $F$  once. This is not a statistical task, as there are no data. Under suitable conditions, the considered zeroth-order method achieves the rate  $d^2/k$ , with  $k$  the number of iterations and it is moreover shown that this rate is optimal. Furthermore, optimal rates are also known in the case that  $f$  is strongly convex and smooth [1, 27].

These rates can further be contrasted to those for stochastic first-order feedback, where the feedback for a query  $\theta$  is the stochastic gradient  $\nabla f(\theta, \Xi)$ . Here [3] show a gap rate of order  $\frac{\sigma^2 d + \|\hat{\theta}_0 - \theta_*\|^2}{k}$  for the average iterate of SGD with constant step size. In [11] these rates are further improved with acceleration to order  $\frac{\sigma^2 d}{k} + \frac{\|\hat{\theta}_0 - \theta_*\|^2}{k^2}$ , matching lower bounds in both contributions. The results were extended beyond the least squares setting in [18] and the related problem of logistic regression was analyzed in the stochastic optimization [4] and bandit settings [15]. For logistic regression the function  $f$  is Lipschitz and higher-order smooth but not strongly convex.

## A Proofs for the upper bound

*Proof of (7):* Expanding the squares yields

$$\boldsymbol{\theta}_k = \boldsymbol{\theta}_{k-1} + \alpha_k \left( 2(Y_k - \mathbf{X}_k^\top \boldsymbol{\theta}_{k-1}) \mathbf{X}_k^\top (\mathbf{U}'_k - \mathbf{U}_k) + (\mathbf{X}_k^\top \mathbf{U}_k)^2 - (\mathbf{X}_k^\top \mathbf{U}'_k)^2 \right) (e^{-\mathbf{U}_k} - e^{\mathbf{U}_k})$$

Setting  $\mathbf{D}_k := e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}$ , the update can be rewritten as affine function in  $\boldsymbol{\theta}_k$ ,

$$\boldsymbol{\theta}_k = \left( I - 2\alpha_k \mathbf{D}_k (\mathbf{U}'_k - \mathbf{U}_k)^\top \mathbf{X}_k \mathbf{X}_k^\top \right) \boldsymbol{\theta}_{k-1} + \alpha_k \left( 2Y_k \mathbf{X}_k^\top (\mathbf{U}'_k - \mathbf{U}_k) + (\mathbf{X}_k^\top \mathbf{U}_k)^2 - (\mathbf{X}_k^\top \mathbf{U}'_k)^2 \right) \mathbf{D}_k,$$

or

$$\begin{aligned}\boldsymbol{\theta}_k - \boldsymbol{\theta}^* &= \left( I - 2\alpha_k \mathbf{D}_k (\mathbf{U}'_k - \mathbf{U}_k)^\top \mathbf{X}_k \mathbf{X}_k^\top \right) (\boldsymbol{\theta}_{k-1} - \boldsymbol{\theta}^*) \\ &\quad + 2\alpha_k (Y_k - \mathbf{X}_k^\top \boldsymbol{\theta}^*) \mathbf{X}_k^\top (\mathbf{U}'_k - \mathbf{U}_k) \mathbf{D}_k \\ &\quad + \alpha_k \left( (\mathbf{X}_k^\top \mathbf{U}_k)^2 - (\mathbf{X}_k^\top \mathbf{U}'_k)^2 \right) \mathbf{D}_k.\end{aligned}$$

Noticing that  $\varepsilon_k = Y_k - \mathbf{X}_k^\top \boldsymbol{\theta}^*$ , (7) follows.  $\square$

*Proof of (8):* We show that  $\mathbb{E}[\boldsymbol{\xi}_k] = 0$ . By definition  $\boldsymbol{\xi}_k := 2\alpha_k \varepsilon_k \mathbf{X}_k^\top (\mathbf{U}'_k - \mathbf{U}_k) (e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}) + \alpha_k \left( (\mathbf{X}_k^\top \mathbf{U}_k)^2 - (\mathbf{X}_k^\top \mathbf{U}'_k)^2 \right) (e^{-\mathbf{U}_k} - e^{\mathbf{U}_k})$ . The first term has expectation zero, since  $\mathbb{E}[\varepsilon_k] = 0$  and  $\varepsilon_k$  is independent of all the other variables. Since  $\mathbf{U}_k$  and  $-\mathbf{U}_k$  have the same distribution,  $\mathbb{E}[e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}] = 0$  and  $\mathbb{E}[(\mathbf{X}_k^\top \mathbf{U}_k)^2 (e^{-\mathbf{U}_k} - e^{\mathbf{U}_k})] = 0$ . Using the independence of  $\mathbf{U}_k, \mathbf{U}'_k$ ,

$$\mathbb{E}[\boldsymbol{\xi}_k] = \mathbb{E}[\alpha_k \left( (\mathbf{X}_k^\top \mathbf{U}_k)^2 - (\mathbf{X}_k^\top \mathbf{U}'_k)^2 \right) (e^{-\mathbf{U}_k} - e^{\mathbf{U}_k})] = 0.$$

## A.1 Expectations with respect to the uniform distribution

The moments of the noise  $U$  play an important role in the analysis.

**Definition 1.** For natural numbers  $r, q \geq 0$ , we abbreviate

$$c_{r,q} := \mathbb{E} [U^r (e^{-U} - e^U)^q] \quad \text{with } U \sim \text{Unif}[-A, A].$$

We have  $c_{0,0} = 1$  and  $c_{2,0} = A^2/3$ . By definition,  $-U(e^{-U} - e^U)$  is a non-negative random variable and therefore  $-c_{1,1} > 0$ . More specifically,

**Lemma 3.** If  $U \sim \text{Unif}([-A, A])$  and  $A \leq 1$ , then

$$\frac{A^2}{11} \leq -c_{1,1} \leq A^2 \quad \text{and} \quad c_{0,2} \leq -3c_{1,1}. \quad (21)$$

*Proof.* Integration by parts gives

$$\begin{aligned}-\mathbb{E}[U(e^{-U} - e^U)] &= -\frac{1}{2A} \int_{-A}^A u(e^{-u} - e^u) du = \frac{1}{A} \int_{-A}^A u e^u du = \frac{1}{A} u e^u \Big|_{-A}^A - \frac{1}{A} \int_{-A}^A e^u du \\ &= e^A + e^{-A} - \frac{e^A - e^{-A}}{A}.\end{aligned}$$

For the second part observe that for  $A \geq 0$ , third order Taylor expansion gives  $e^A \geq 1 + A + A^2/2$  and  $e^{-A} \leq 1 - A + A^2/2$ . From the expression above and using that  $A \leq 1$ ,

$$\begin{aligned}-\mathbb{E}[U(e^{-U} - e^U)] &= e^A + e^{-A} - \frac{e^A - e^{-A}}{A} \\ &= -e^A \left( \frac{1}{A} - 1 \right) + e^{-A} \left( 1 + \frac{1}{A} \right) \\ &\leq -\left( 1 + A + \frac{A^2}{2} \right) \frac{1-A}{A} + \left( 1 - A + \frac{A^2}{2} \right) \frac{A+1}{A} \\ &= -(1+A) \frac{1-A}{A} - \frac{A-A^2}{2} + (1-A) \frac{A+1}{A} + \frac{A^2+A}{2} \\ &= A^2.\end{aligned}$$

By third order Taylor expansion and  $A \leq 1$ , we find  $e^A \leq 1 + A + A^2/2 + e^A A^3/6 \leq 1 + A + A^2/2 + eA^3/6$  and  $e^{-A} \geq 1 - A + A^2/2 - eA^3/6$ . Basically following the same steps as for the upper bound of  $-\mathbb{E}[U(e^{-U} - e^U)]$ ,

$$\begin{aligned} -\mathbb{E}[U(e^{-U} - e^U)] &= -e^A \left( \frac{1}{A} - 1 \right) + e^{-A} \left( 1 + \frac{1}{A} \right) \\ &\geq - \left( 1 + A + \frac{A^2}{2} + e \frac{A^3}{6} \right) \frac{1-A}{A} + \left( 1 - A + \frac{A^2}{2} - e \frac{A^3}{6} \right) \frac{A+1}{A} \\ &= A^2 \left( 1 - \frac{e}{3} \right). \end{aligned}$$

Since  $1 - e/3 \geq 0.093 \geq 1/11$ , this completes the proof for (21).

To prove  $c_{0,2} \leq -3c_{1,1}$ , we use that  $e^x - e^{-x} = 2 \sum_{\ell \text{ odd}} x^\ell / \ell!$ . Since for odd  $\ell$ ,  $x^\ell (e^x - e^{-x}) \geq 0$ , we have for  $|x| \leq 1$  that  $x^\ell (e^x - e^{-x}) \leq x(e^x - e^{-x})$ . With  $e \leq 3$ , we obtain for  $x \leq 1$ ,

$$\begin{aligned} (e^x - e^{-x})^2 &= 2 \sum_{\ell \text{ odd}} \frac{x^\ell}{\ell!} (e^x - e^{-x}) \leq 2 \sum_{\ell \text{ odd}} \frac{1}{\ell!} x (e^x - e^{-x}) \\ &\leq 2 \left( e - 1 - \frac{1}{2!} \right) x (e^x - e^{-x}) \leq 3x (e^x - e^{-x}). \end{aligned}$$

Thus  $c_{0,2} = \mathbb{E}[(e^U - e^{-U})^2] \leq 3\mathbb{E}[U(e^U - e^{-U})] = -3c_{1,1}$ .  $\square$

Using the previous lemma,

$$c_{r,2} \leq A^r c_{0,2} \leq 3A^{r+2}. \quad (22)$$

For odd  $r$ ,  $c_{r,2} = 0$  and this inequality makes only sense if  $r$  is even.

## A.2 Moments of the noise contributions in the zeroth-order scheme

Recall that  $\|\cdot\|$  denotes the spectral norm for matrices and the Euclidean norm for vectors.

**Lemma 4.** *Let  $\mathbf{a}, \mathbf{b}$  be column vectors of the same length. Then  $\|\mathbf{a}\mathbf{b}^\top\| = \|\mathbf{a}\| \|\mathbf{b}\|$ .*

*Proof.* The matrix  $\mathbf{b}\mathbf{b}^\top$  is of rank one with non-zero eigenvalue  $\mathbf{b}^\top \mathbf{b}$  and corresponding eigenvector  $\|\mathbf{b}\|^{-1} \mathbf{b}$ . Thus,  $\|\mathbf{a}\mathbf{b}^\top\|^2 = \lambda_{\max}(\mathbf{b}\mathbf{a}^\top \mathbf{a}\mathbf{b}^\top) = \mathbf{a}^\top \mathbf{a} \lambda_{\max}(\mathbf{b}\mathbf{b}^\top) = \mathbf{a}^\top \mathbf{a} \mathbf{b}^\top \mathbf{b}$ . The result follows by taking square roots.  $\square$

We now derive closed-form expressions and bounds for the two expected values in (10). Recall that  $c_{r,q} = \mathbb{E}[U^r (e^{-U} - e^U)^q]$ . The  $k$ -th power  $\mathbf{x}^k$  of a vector  $\mathbf{x} = (x_1, \dots, x_d)$  is understood componentwise, e.g.,  $\|\mathbf{x}^2\|^2 = \sum_i x_i^4$ .

**Lemma 5.** *Let  $I$  denote the  $d \times d$  identity matrix. We have*

$$V(\mathbf{x}) = (c_{2,2} - c_{2,0}c_{0,2} - 2c_{1,1}^2) \text{diag}(\mathbf{x}^2) + 2c_{1,1}^2 \mathbf{x}\mathbf{x}^\top + 2c_{2,0}c_{0,2} \|\mathbf{x}\|^2 I.$$

and

$$\begin{aligned}
W &= \mathbb{E} \left[ (c_{4,2} - 6c_{2,0}c_{2,2} + 6c_{2,0}^2c_{0,2} - c_{4,0}c_{0,2} - 8c_{3,1}c_{1,1} + 24c_{2,0}c_{1,1}^2) \text{diag}(\mathbf{X}^4) \right. \\
&\quad + (4c_{2,0}c_{2,2} - 4c_{2,0}^2c_{0,2} - 8c_{2,0}c_{1,1}^2) \|\mathbf{X}\|^2 \text{diag}(\mathbf{X}^2) \\
&\quad + (4c_{3,1}c_{1,1} - 12c_{2,0}c_{1,1}^2)(\mathbf{X}^3\mathbf{X}^\top + \mathbf{X}(\mathbf{X}^3)^\top) \\
&\quad + 8c_{2,0}c_{1,1}^2 \|\mathbf{X}\|^2 \mathbf{X}\mathbf{X}^\top \\
&\quad + (2c_{4,0}c_{0,2} - 6c_{2,0}^2c_{0,2}) \|\mathbf{X}^2\|^2 I \\
&\quad \left. + 4c_{2,0}^2c_{0,2} \|\mathbf{X}\|^4 I \right].
\end{aligned}$$

Moreover for  $\mathbf{X} = (X_1, \dots, X_d)^\top$ ,

$$\mathbb{E}[V(\mathbf{X})] \leq 12A^4d \max_{i=1, \dots, d} \mathbb{E}[X_i^2]I, \quad (23)$$

$$\mathbb{E}[\|\mathbf{X}\|^2 V(\mathbf{X})] \leq 12A^4d^2 \max_{i=1, \dots, d} \mathbb{E}[X_i^4]I, \quad (24)$$

$$W \leq 107A^6d^2 \max_{i=1, \dots, d} 1 \vee \mathbb{E}[X_i^4]I. \quad (25)$$

While the closed-form expression of  $W$  depends on fourth power of  $\mathbf{X}$ , it is convenient to relate this to a sixth power in (25).

*Proof.* We first prove the formula for  $W$ . Let  $\mathbf{x} = (x_1, \dots, x_d)$  be fixed. Using that  $\mathbb{E}[(\mathbf{x}^\top \mathbf{U}')^2] = \mathbb{E}[(\mathbf{U}')^\top \mathbf{x}\mathbf{x}^\top \mathbf{U}'] = c_{2,0} \text{tr}(\mathbf{x}\mathbf{x}^\top) = c_{2,0}\mathbf{x}^\top \mathbf{x}$ ,

$$\begin{aligned}
\mathbb{E}[(\mathbf{x}^\top \mathbf{U}')^4] &= \mathbb{E} \left[ \sum_{i,j,k,\ell} U'_i U'_j U'_k U'_\ell x_i x_j x_k x_\ell \right] \\
&= c_{4,0} \sum_i x_i^4 + 3c_{2,0}^2 \sum_{i \neq j} x_i^2 x_j^2 \\
&= (c_{4,0} - 3c_{2,0}^2) \|\mathbf{x}^2\|^2 + 3c_{2,0}^2 \|\mathbf{x}\|^4,
\end{aligned}$$

$\mathbb{E}[\mathbf{D}\mathbf{D}^\top] = c_{0,2}I$ , and the independence of  $\mathbf{U}$  and  $\mathbf{U}'$  gives

$$\begin{aligned}
&\mathbb{E} \left[ ((\mathbf{x}^\top \mathbf{U})^2 - (\mathbf{x}^\top \mathbf{U}')^2)^2 \mathbf{D}\mathbf{D}^\top \right] \\
&= \mathbb{E} \left[ \left( (\mathbf{x}^\top \mathbf{U})^4 - 2c_{2,0}(\mathbf{x}^\top \mathbf{U})^2 \mathbf{x}^\top \mathbf{x} + (c_{4,0} - 3c_{2,0}^2) \|\mathbf{x}^2\|^2 + 3c_{2,0}^2 \|\mathbf{x}\|^4 \right) \mathbf{D}\mathbf{D}^\top \right] \\
&= \mathbb{E} [(\mathbf{x}^\top \mathbf{U})^4 \mathbf{D}\mathbf{D}^\top] - 2c_{2,0} \mathbf{x}^\top \mathbf{x} \mathbb{E} [(\mathbf{x}^\top \mathbf{U})^2 \mathbf{D}\mathbf{D}^\top] + c_{0,2} \left( (c_{4,0} - 3c_{2,0}^2) \|\mathbf{x}^2\|^2 + 3c_{2,0}^2 \|\mathbf{x}\|^4 \right) I.
\end{aligned}$$

Next we simplify these two remaining expectations. The  $(i, j)$ -th off-diagonal entry of the matrix  $\mathbb{E}[(\mathbf{x}^\top \mathbf{U})^2 \mathbf{D}\mathbf{D}^\top]$  is  $\mathbb{E}[\sum_{\ell,k} U_\ell U_k x_\ell x_k D_i D_j]$ . The summands are non-zero if either  $(\ell, k) = (i, j)$  or  $(\ell, k) = (j, i)$ . In both cases we get the contribution  $c_{1,1}^2 x_i x_j$ , such that for  $i \neq j$ ,  $\mathbb{E}[\sum_{\ell,k} U_\ell U_k x_\ell x_k D_i D_j] = 2c_{1,1}^2 x_i x_j$ . For the  $i$ -th diagonal entry of  $\mathbb{E}[(\mathbf{x}^\top \mathbf{U})^2 \mathbf{D}\mathbf{D}^\top]$ , we obtain  $\mathbb{E}[\sum_{\ell,k} U_\ell U_k x_\ell x_k D_i^2] = \mathbb{E}[\sum_{\ell} U_\ell^2 x_\ell^2 D_i^2] = c_{2,2} x_i^2 + c_{2,0} c_{0,2} \sum_{\ell \neq i} x_\ell^2 = (c_{2,2} - c_{2,0} c_{0,2}) x_i^2 + c_{2,0} c_{0,2} \|\mathbf{x}\|^2$ . Combining these formulas yields

$$\mathbb{E}[(\mathbf{x}^\top \mathbf{U})^2 \mathbf{D}\mathbf{D}^\top] = (c_{2,2} - c_{2,0} c_{0,2} - 2c_{1,1}^2) \text{diag}(\mathbf{x}^2) + 2c_{1,1}^2 \mathbf{x}\mathbf{x}^\top + c_{2,0} c_{0,2} \|\mathbf{x}\|_2^2 I \quad (26)$$



and

$$\begin{aligned}
\mathbb{E}[(\mathbf{x}^\top \mathbf{U})^4 \mathbf{D} \mathbf{D}^\top] &= \mathbb{E} \left[ \sum_{i,j,\ell,m} x_i U_i x_j U_j x_\ell U_\ell x_m U_m \mathbf{D} \mathbf{D}^\top \right] \\
&= (c_{4,2} - 6c_{2,0}c_{2,2} + 6c_{2,0}^2 c_{0,2} - c_{4,0}c_{0,2} - 8c_{3,1}c_{1,1} + 24c_{2,0}c_{1,1}^2) \text{diag}(\mathbf{x}^4) \\
&\quad + (6c_{2,0}c_{2,2} - 6c_{2,0}^2 c_{0,2} - 12c_{2,0}c_{1,1}^2) \|\mathbf{x}\|^2 \text{diag}(\mathbf{x}^2) \\
&\quad + (c_{4,0}c_{0,2} - 3c_{2,0}^2 c_{0,2}) \|\mathbf{x}^2\|^2 I \\
&\quad + 3c_{2,0}^2 c_{0,2} \|\mathbf{x}\|^4 I \\
&\quad + 12c_{2,0}c_{1,1}^2 \|\mathbf{x}\|^2 \mathbf{x} \mathbf{x}^\top \\
&\quad + (4c_{3,1}c_{1,1} - 12c_{2,0}c_{1,1}^2) (\mathbf{x}^3 \mathbf{x}^\top + \mathbf{x}(\mathbf{x}^3)^\top).
\end{aligned}$$

This is because on a diagonal entry  $(p, p)$ , we have  $c_{4,2}x_p^4$ , as well as  $6c_{2,0}c_{2,2}x_p^2 \sum_{q \neq p} x_q^2$  and  $3c_{2,0}^2 c_{0,2} (\sum_{q \neq p} x_q^2)^2$  and  $c_{4,0}c_{0,2} \sum_{q \neq p} x_q^4$ . On an off-diagonal entry  $(p, q)$ , we have  $12c_{2,0}c_{1,1}^2 x_p x_q \sum_{r \notin \{p, q\}} x_r^2$  as well as  $4c_{3,1}c_{1,1} (x_p^3 x_q + x_p x_q^3)$ . Combining/grouping terms, replacing  $\mathbf{x}$  by the random vector  $\mathbf{X}$  and taking expectation with respect to  $\mathbf{X}$  yields the formula for  $W$ .

For  $V(\mathbf{x})$ , observe that the cross-terms are zero and by conditioning first on  $\mathbf{U}'$  and using that  $\mathbb{E}[\mathbf{D} \mathbf{D}^\top] = c_{0,2} I$ , we find

$$\begin{aligned}
V(\mathbf{x}) &= \mathbb{E}[\mathbf{x}^\top (\mathbf{U}' - \mathbf{U}) \mathbf{D} \mathbf{D}^\top (\mathbf{U}' - \mathbf{U})^\top \mathbf{x}] \\
&= \mathbb{E}[\mathbf{x}^\top \mathbf{U}' \mathbf{D} \mathbf{D}^\top (\mathbf{U}')^\top \mathbf{x}] + \mathbb{E}[\mathbf{x}^\top \mathbf{U} \mathbf{D} \mathbf{D}^\top \mathbf{U}^\top \mathbf{x}] \\
&= c_{0,2} \mathbf{x}^\top \mathbb{E}[\mathbf{U}' (\mathbf{U}')^\top] \mathbf{x} + \mathbb{E}[\mathbf{x}^\top \mathbf{U} \mathbf{D} \mathbf{D}^\top \mathbf{U}^\top \mathbf{x}] \\
&= c_{0,2} c_{2,0} \|\mathbf{x}\|^2 I + \mathbb{E}[\mathbf{x}^\top \mathbf{U} \mathbf{D} \mathbf{D}^\top \mathbf{U}^\top \mathbf{x}].
\end{aligned}$$

Combined with (26), the formula  $V(\mathbf{x}) = (c_{2,2} - c_{2,0}c_{0,2} - 2c_{1,1}^2) \text{diag}(\mathbf{x}^2) + 2c_{1,1}^2 \mathbf{x} \mathbf{x}^\top + 2c_{2,0}c_{0,2} \|\mathbf{x}\|^2 I$  follows.

For the bounds on the expectations, we use the bounds on the moments  $c_{r,q}$  derived in Section A.1. In particular  $c_{r,q} \geq 0$  whenever  $r$  and  $q$  are even. The matrix  $\mathbf{x} \mathbf{x}^\top$  is positive semi-definite. By (21),  $c_{1,1} \leq A^2$  and  $c_{0,2} \leq 3A^2$ . Thus,  $V(\mathbf{x}) \leq 3A^4 \text{diag}(\mathbf{x}^2) + 2A^4 \mathbf{x} \mathbf{x}^\top + 6A^4 \|\mathbf{x}\|^2 I$ . Moreover, by Lemma 4, the largest eigenvalue of  $\mathbf{x} \mathbf{x}^\top$  is  $\|\mathbf{x}\|^2$  and therefore  $\mathbf{x} \mathbf{x}^\top \leq \|\mathbf{x}\|^2 I$ . Since also  $\text{diag}(\mathbf{x}^2) \leq \|\mathbf{x}\|^2 I$ , we obtain  $V(\mathbf{x}) \leq 12A^4 \|\mathbf{x}\|^2 I$  and thus for  $\mathbf{X} = (X_1, \dots, X_d)^\top$ ,  $\mathbb{E}[V(\mathbf{X})] \leq 12A^4 d \max_{i=1, \dots, d} \mathbb{E}[X_i^2] I$ , proving (23). Since  $X_i^2 X_j^2 \leq X_i^4/2 + X_j^4/2$ , we can derive moreover  $\mathbb{E}[\|\mathbf{X}\|^2 V(\mathbf{X})] \leq 12A^4 d^2 \max_{i,j=1, \dots, d} \mathbb{E}[X_i^2 X_j^2] \leq 12A^4 d^2 \max_{i=1, \dots, d} \mathbb{E}[X_i^4]$ , proving (24).

We finally derive (25). Since  $2\|\mathbf{x}^3 \mathbf{x}^\top\| = 2\|\mathbf{x}^3\| \|\mathbf{x}\|$ , all eigenvalues of  $\mathbf{x}^3 \mathbf{x}^\top + \mathbf{x}(\mathbf{x}^3)^\top$  lie between  $-2\|\mathbf{x}^3\| \|\mathbf{x}\|$  and  $2\|\mathbf{x}^3\| \|\mathbf{x}\|$ . Therefore, for any real number  $a$ , we have  $a(\mathbf{x}^3 \mathbf{x}^\top + \mathbf{x}(\mathbf{x}^3)^\top) \leq 2|a| \|\mathbf{x}^3\| \|\mathbf{x}\|$ . By (22),  $c_{r,2} \leq A^r c_{0,2} \leq 3A^{r+2}$ . Moreover  $0 \leq c_{3,1}c_{1,1} \leq A^2 c_{1,1}^2 \leq A^6$  and  $0 \leq c_{2,0}c_{1,1}^2 \leq A^6$ . Thus,

$$\begin{aligned}
W &\leq \mathbb{E} \left[ (3 + 18 + 24)A^6 \text{diag}(\mathbf{X}^4) + 12A^6 \|\mathbf{X}\|^2 \text{diag}(\mathbf{X}^2) + 12A^6 2\|\mathbf{X}^3\| \|\mathbf{X}\| I + 8A^6 \|\mathbf{X}\|^4 I \right. \\
&\quad \left. + 6A^6 \|\mathbf{X}^2\|^2 I + 12A^6 \|\mathbf{X}\|^4 I \right].
\end{aligned}$$

Since  $\text{diag}(\mathbf{X}^4) \leq \|\mathbf{X}\|^4 I$ ,  $\text{diag}(\mathbf{X}^2) \leq \|\mathbf{X}\|^2 I$ , and  $\|\mathbf{X}^2\|^2 \leq \|\mathbf{X}\|^4$ , we obtain  $W \leq 83A^6 \mathbb{E}[\|\mathbf{X}\|^4] I + 24A^6 \mathbb{E}[\|\mathbf{X}^3\| \|\mathbf{X}\|] I$ .

Because of the elementary inequality  $ab \leq \frac{1}{2}a^2 + \frac{1}{2}b^2$ ,

$$\mathbb{E}[\|\mathbf{X}\|^4] = \sum_{i,j=1}^d \mathbb{E}[X_i^2 X_j^2] \leq \frac{1}{2} \sum_{i,j=1}^d \mathbb{E}[X_i^4] + \mathbb{E}[X_j^4] \leq d^2 \max_{i=1,\dots,d} \mathbb{E}[X_i^4].$$

Similarly,  $\|\mathbf{X}^3\| \|\mathbf{X}\| \leq \sqrt{d \max_i X_i^6} \sqrt{d \max_i X_i^2} = d \max_i X_i^4 \leq d \sum_i X_i^4$  and hence

$$\mathbb{E}[\|\mathbf{X}^3\| \|\mathbf{X}\|] \leq d \sum_i \mathbb{E}[X_i^4] \leq d^2 \max_i \mathbb{E}[X_i^4].$$

Combined with  $W \leq 83A^6 \mathbb{E}[\|\mathbf{X}\|^4] I + 24A^6 \mathbb{E}[\|\mathbf{X}^3\| \|\mathbf{X}\|] I$ , we finally obtain

$$W \leq 107A^6 d^2 \max_{i=1,\dots,d} \mathbb{E}[X_i^4] I.$$

□

### A.3 Proof of Lemma 1

Recall that  $\mathbf{D}_k := e^{-\mathbf{U}_k} - e^{\mathbf{U}_k}$ . The rewritten update equation (7) can be decomposed into the following three parts

$$\boldsymbol{\theta}_k - \boldsymbol{\theta}^* = \mathbf{A} + \mathbf{B} + \mathbf{C}, \quad (27)$$

where

$$\begin{aligned} \mathbf{A} &:= (I - 2\alpha_k \mathbf{D}_k (\mathbf{U}'_k - \mathbf{U}_k)^\top \mathbf{X}_k \mathbf{X}_k^\top) (\boldsymbol{\theta}_{k-1} - \boldsymbol{\theta}^*), \\ \mathbf{B} &:= 2\alpha_k \varepsilon_k \mathbf{X}_k^\top (\mathbf{U}'_k - \mathbf{U}_k) \mathbf{D}_k, \\ \mathbf{C} &:= \alpha_k ((\mathbf{X}_k^\top \mathbf{U}_k)^2 - (\mathbf{X}_k^\top \mathbf{U}'_k)^2) \mathbf{D}_k. \end{aligned}$$

In the following we often use that  $(\mathbf{X}_k, Y_k)$  has the same distribution as  $(\mathbf{X}, Y)$ . This means that if all randomness in an expectation is only due to  $(\mathbf{X}_k, Y_k)$ , these variables can be replaced by  $(\mathbf{X}, Y)$ . Recall moreover that odd powers of  $\mathbf{U}$  or  $\mathbf{U}'$  disappear and that we have defined  $\mu = -c_{1,1}$ . Thus,

$$\begin{aligned} \mathbb{E}[\mathbf{A}\mathbf{A}^\top] &= S_{k-1} - 2\mu\alpha_k(QS_{k-1} + S_{k-1}Q) + 4\alpha_k^2 \mathbb{E}[\mathbf{X}_k^\top S_{k-1} \mathbf{X}_k V(\mathbf{X}_k)] \\ &= (I - 2\alpha_k \mu Q) S_{k-1} (I - 2\alpha_k \mu Q) + 4\alpha_k^2 \{ \mathbb{E}[\mathbf{X}^\top S_{k-1} \mathbf{X} V(\mathbf{X})] - \mu^2 Q S_{k-1} Q \}. \end{aligned} \quad (28)$$

Since for any  $d \times d$  matrix  $T$ , we have  $T^\top S_{k-1} T \leq \|S_{k-1}\| T^\top T$ , we also have

$$\mathbb{E}[\mathbf{A}\mathbf{A}^\top] \leq \|S_{k-1}\| (I - 4\alpha_k \mu Q + 4\alpha_k^2 \mathbb{E}[\|\mathbf{X}\|^2 V(\mathbf{X})]). \quad (29)$$

Now by counting powers of  $\mathbf{U}_k, \mathbf{U}'_k, \mathbf{D}_k$ , we find that

$$\mathbb{E}[\mathbf{A}\mathbf{C}^\top] = 0.$$

Applying tower rule by conditioning on all randomness except  $\varepsilon_k \sim \mathcal{N}(0, \sigma^2)$ , it follows that

$$\begin{aligned}\mathbb{E}[\mathbf{AB}^\top] &= \mathbb{E}[\mathbb{E}[\mathbf{AB}^\top \mid \mathbf{U}_k, \mathbf{U}'_k, \mathbf{X}_k, \boldsymbol{\theta}_k]] = \mathbb{E}[\mathbf{A}2\alpha_k\mathbf{D}_k^\top(\mathbf{U}'_k - \mathbf{U}_k)^\top \mathbf{X}_k \mathbb{E}[\varepsilon_k]] = 0, \\ \mathbb{E}[\mathbf{BC}^\top] &= 0, \\ \mathbb{E}[\mathbf{BB}^\top] &= 4\alpha_k^2 \mathbb{E}[\varepsilon_k^2 V(\mathbf{X})] = 4\alpha_k^2 \sigma^2 \mathbb{E}[V(\mathbf{X})],\end{aligned}$$

and finally

$$\mathbb{E}[\mathbf{CC}^\top] = \alpha_k^2 W,$$

which is some combination of fourth powers of  $\mathbf{X}$ .

Now (i) follows due to  $\mathbb{E}[G_k \mathbf{W}_k \boldsymbol{\xi}_k^\top] = \mathbb{E}[\mathbf{A}(\mathbf{B} + \mathbf{C})^\top] = 0$ . Statement (ii) follows since by (8),  $\mathbb{E}[\boldsymbol{\xi}_k] = 0$  such that  $\text{Cov}(\boldsymbol{\xi}_k) = \mathbb{E}[\boldsymbol{\xi}_k \boldsymbol{\xi}_k^\top] = \mathbb{E}[(\mathbf{B} + \mathbf{C})(\mathbf{B} + \mathbf{C})^\top] = \mathbb{E}[\mathbf{BB}^\top] + \mathbb{E}[\mathbf{CC}^\top] = 4\alpha_k^2 \sigma^2 \mathbb{E}[V(\mathbf{X})] + \alpha_k^2 W$ . To see (iii), observe that

$$\begin{aligned}S_k &= \mathbb{E}[(\mathbf{A} + \mathbf{B} + \mathbf{C})(\mathbf{A} + \mathbf{B} + \mathbf{C})^\top] \\ &= \mathbb{E}[\mathbf{AA}^\top] + \mathbb{E}[\mathbf{AB}^\top] + \mathbb{E}[\mathbf{BA}^\top] + \mathbb{E}[\mathbf{BB}^\top] + \mathbb{E}[\mathbf{CC}^\top] \\ &= (I - 2\alpha_k \mu Q) S_{k-1} (I - 2\alpha_k \mu Q) + 4\alpha_k^2 \left( \mathbb{E}[\mathbf{X}^\top S_{k-1} \mathbf{X} V(\mathbf{X})] - \mu^2 Q S_{k-1} Q \right) \\ &\quad + 4\alpha_k^2 \sigma^2 \mathbb{E}[V(\mathbf{X})] + \alpha_k^2 W.\end{aligned}$$

Applying (29) instead of (28) yields the asserted inequality.  $\square$

#### A.4 Proof of Theorem 1

We now combine the recursive formula with the bounds obtained in (23), (24), and (25). Recall that by definition  $\mu = -c_{1,1}$  and by Lemma 3,  $\mu \geq A^2/11$  which is the same as  $-\mu \leq -A^2/11$ . Setting  $M_k := \max_{i=1, \dots, d} 1 \vee E[X_i^k]$ , we obtain

$$\|S_k\| \leq \|S_{k-1}\| \left( 1 - \frac{4}{11} \alpha_k A^2 \lambda_{\min}(Q) + 48\alpha_k^2 A^4 d^2 M_4 \right) + 48\alpha_k^2 \sigma^2 A^4 d M_2 + \alpha_k^2 107 A^6 d^2 M_4.$$

By assumption  $A = \sigma/\sqrt{d}$ . If

$$96\alpha_k^2 A^4 d^2 M_4 \leq \frac{4}{11} \alpha_k A^2 \lambda_{\min}(Q) \leq 1, \quad (30)$$

then, the previous inequality yields the simpler recursion inequality

$$\begin{aligned}\|S_k\| &\leq \|S_{k-1}\| \left( 1 - \frac{2}{11} \alpha_k A^2 \lambda_{\min}(Q) \right) + (48M_2 + 107M_4) \alpha_k^2 \sigma^2 A^4 d \\ &= \|S_{k-1}\| (1 - \beta_k) + C_* \beta_k^2 \sigma^2 d,\end{aligned}$$

where we defined

$$\beta_k := \frac{2}{11} \alpha_k A^2 \lambda_{\min}(Q) \quad \text{and} \quad C_* := (48M_2 + 107M_4) \left( \frac{11}{2\lambda_{\min}(Q)} \right)^2. \quad (31)$$

Induction with respect to  $k$  gives

$$\|S_k\| \leq \|S_0\| \prod_{\ell=1}^k (1 - \beta_\ell) + C_* \sum_{j=1}^k \left( \prod_{\ell=j+1}^k (1 - \beta_\ell) \right) \beta_j^2 \sigma^2 d, \quad (32)$$

where  $\prod_{\ell=k+1}^k (1 - \beta_\ell) := 1$ . Rewritten in terms of  $\beta_k$ , the assumed inequality (30) is

$$96 \left( \frac{11\beta_k^2}{2\lambda_{\min}(Q)} \right)^2 d^2 M_4 \leq 2\beta_k \leq 1.$$

The choice of the learning rate for  $\alpha_k$  in (12) results in

$$\beta_k = \frac{2}{11} \alpha_k A^2 \lambda_{\min}(Q) = \frac{2B \log(d)}{Bk + d^2 \log(d)} \quad \text{with } B = 1 \wedge \frac{\lambda_{\min}(Q)^2}{2904M_4}. \quad (33)$$

For all  $d \geq 1$ , we have  $\beta_k \leq 2B/d^2$ . Using that  $96(11/2)^2 = 2904$ , one can now check that for all  $d \geq 2$ , the assumed inequalities for  $\beta_k$  hold and afortiori thus also (30).

Let  $k^*$  be the smallest integer such that  $Bk^* \geq d^2 \log(d)$ . Since  $d \geq 2$ , we must have

$$k^* \leq \frac{1 + d^2 \log(d)}{B} \leq \frac{2d^2 \log(d)}{B}. \quad (34)$$

For all  $k \geq k^*$ , we have  $\beta_k \geq \log(d)/k$ . By bounding  $1 - \beta_\ell \leq 1$  for all  $\ell < k^*$ , and using for  $k \geq j + 1 \geq k^*$ ,  $\sum_{\ell=j+1}^k \frac{1}{\ell} \geq \sum_{\ell=j+1}^k \int_{\ell}^{\ell+1} \frac{1}{u} du = \int_{j+1}^{k+1} \frac{1}{u} du = \log(k+1) - \log(j+1) = \log((k+1)/(j+1))$ , we obtain for all  $k \geq k^*$ ,

$$\prod_{\ell=j+1}^k (1 - \beta_\ell) \leq \exp \left( -\log(d) \sum_{\ell=k^* \vee (j+1)}^k \frac{1}{\ell} \right) \leq \left( \frac{k^* \vee (j+1)}{k+1} \right)^{\log(d)}.$$

Combined with  $\beta_k \leq 2B/d^2$  for all  $k$ ,  $\beta_j \leq \log(d)/j \leq 2\log(d)/(j+1)$  for all  $j \geq k^*$ , and (32), we obtain for all  $k > k^*$ ,

$$\begin{aligned} \|S_k\| &\leq \left( \frac{k^*}{k} \right)^{\log(d)} \|S_0\| + C_* \sum_{j=1}^{k^*-1} \left( \frac{k^*}{k} \right)^{\log(d)} \left( \frac{2B}{d^2} \right)^2 \sigma^2 d + C_* \sum_{j=k^*}^k \left( \frac{j+1}{k+1} \right)^{\log(d)} \left( \frac{2\log(d)}{j+1} \right)^2 \sigma^2 d \\ &\leq \left( \frac{k^*}{k} \right)^{\log(d)} \|S_0\| + C_* \frac{1}{k} \left( \frac{k^*}{k} \right)^{\log(d)-1} (k^*)^2 \left( \frac{2B}{d^2} \right)^2 \sigma^2 d + C_* \sum_{j=k^*}^k \left( \frac{j+1}{k+1} \right)^{\log(d)} \left( \frac{2\log(d)}{j+1} \right)^2 \sigma^2 d. \end{aligned}$$

Since  $d \geq 9$ , we have  $\log(d) - 2 > 0$  and therefore  $k^*/k \leq 1$ . Using also (34) gives

$$\frac{1}{k} \left( \frac{k^*}{k} \right)^{\log(d)-1} (k^*)^2 \left( \frac{2B}{d^2} \right)^2 \sigma^2 d \leq \frac{1}{k} \left( \frac{2d^2 \log(d)}{B} \right)^2 \left( \frac{2B}{d^2} \right)^2 \sigma^2 d \leq \frac{16\sigma^2 d \log^2(d)}{k}.$$

Applying  $\log(d) - 2 > 0$  again and moreover,  $(j+1)/(k+1) \leq 1$  for all  $j \leq k$ , we find

$$\sum_{j=k^*}^k \left( \frac{j+1}{k+1} \right)^{\log(d)} \frac{1}{(j+1)^2} = \frac{1}{(k+1)^2} \sum_{j=k^*}^k \left( \frac{j+1}{k+1} \right)^{\log(d)-2} \leq \frac{1}{k}.$$

Combined, the last three displayed inequalities combined yield

$$\|S_k\| \leq \left( \frac{k^*}{k} \right)^{\log(d)} \|S_0\| + C_* \frac{20\sigma^2 d \log^2(d)}{k}.$$

The first claim follows now with  $C = 20C_*$ .

For two positive semi-definite matrices  $D = D_1^\top D_1$  and  $E$ , we know that  $D_1 E D_1^\top$  is positive semi-definite and bounded by  $\|E\| D_1 D_1^\top$ . Therefore,  $\text{tr}(DE) = \text{tr}(D_1 E D_1^\top) \leq \text{tr}(\|E\| D_1 D_1^\top) = \|E\| \text{tr}(D_1^\top D_1) = \|E\| \text{tr}(D)$ . Applying this to  $(D, E) = (Q, S_k)$  and using that

$$\mathbb{E}[(\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)^\top Q (\boldsymbol{\theta}_k - \boldsymbol{\theta}^*)] = \text{tr}(Q S_k)$$

as well as  $\text{tr}(Q) \leq d$  yields the second claimed inequality.  $\square$

*Proof of Inequality (14).* Set  $k_* := 2e^{\gamma+3\kappa} d^2 \log(d)/B$ . Then,

$$\left(\frac{2d^2 \log(d)}{B k_*}\right)^{\log(d)} \leq d^{-\gamma-3\kappa} \leq \frac{1}{d^\gamma (d^2 \log(d))^\kappa} \leq \frac{C'}{d^\gamma k_*^\kappa}, \quad (35)$$

with  $C' = (2e^{\gamma+3\kappa}/B)^\kappa$ . If for positive  $A, B, u, v, x$  with  $u > v$ , we have  $A/x^u \leq B/x^v$ , then also  $A/y^u \leq B/y^v$  for all  $y \geq x$ . Since  $d > e^\kappa$  we can apply this inequality with  $u = \log(d) > \kappa = v$ , proving that (35) holds for all  $k \geq k_*$ .

## B Proofs for the lower bounds

*Proof of Lemma 2.* To distinguish the two models, we rename the queries in (16) into  $W_1, W'_1, \dots, W_k, W'_k$  and the query vectors into  $\mathbf{w}_1, \mathbf{w}'_1, \dots, \mathbf{w}_k, \mathbf{w}'_k$ . This means that the second query model is then denoted by  $\{W_1, W'_1, \dots, W_k, W'_k, \mathbf{w}_1, \mathbf{w}'_1, \dots, \mathbf{w}_k, \mathbf{w}'_k\}$  with  $W_\ell = Y_\ell - \mathbf{X}_\ell^\top \mathbf{w}_\ell$  and  $W'_\ell = \mathbf{X}_\ell^\top \mathbf{w}'_\ell$ .

To see the equivalence, we apply induction with respect to  $k$ . The base case  $k = 1$  and the induction step  $k \rightarrow k + 1$  are similar and therefore only the latter will be discussed.

The induction step  $k \rightarrow k + 1$  is split in two parts. We first prove that the first query model can be transformed into the second query model without knowledge of the parameters. To see this, choose query vectors  $\mathbf{v}_{k+1} = \mathbf{w}_{k+1} + \mathbf{w}'_{k+1}$  and  $\mathbf{v}'_{k+1} = \mathbf{w}_{k+1} - \mathbf{w}'_{k+1}$ . Query vectors can depend on previously seen queries and query vectors, such that  $\mathbf{w}_{k+1}, \mathbf{w}'_{k+1}$  can depend on  $\{W_1, W'_1, \dots, W_k, W'_k, \mathbf{w}_1, \mathbf{w}'_1, \dots, \mathbf{w}_k, \mathbf{w}'_k\}$  and  $\mathbf{v}_{k+1}, \mathbf{v}'_{k+1}$  can depend on  $\{Z_1, Z'_1, \dots, Z_k, Z'_k, \mathbf{v}_1, \mathbf{v}'_1, \dots, \mathbf{v}_k, \mathbf{v}'_k\}$ . Since the models can be transformed into each other by the induction hypothesis, those are eligible query vectors. The corresponding queries in (15) are then given by  $Z_{k+1} = Y_{k+1} - \mathbf{X}_{k+1}^\top (\mathbf{w}_{k+1} + \mathbf{w}'_{k+1})$  and  $Z'_{k+1} = Y_{k+1} - \mathbf{X}_{k+1}^\top (\mathbf{w}_{k+1} - \mathbf{w}'_{k+1})$ . Now,  $W_{k+1} := \frac{1}{2}(Z_{k+1} + Z'_{k+1}) = Y_{k+1} - \mathbf{X}_{k+1}^\top \mathbf{w}_{k+1}$  and  $W'_{k+1} := \frac{1}{2}(Z'_{k+1} - Z_{k+1}) = \mathbf{X}_{k+1}^\top \mathbf{w}'_{k+1}$  are then the queries from the second model. Since one can also retrieve the query vectors  $\mathbf{w}_{k+1} = \frac{1}{2}(\mathbf{v}_{k+1} + \mathbf{v}'_{k+1})$  and  $\mathbf{w}'_{k+1} = \frac{1}{2}(\mathbf{v}_{k+1} - \mathbf{v}'_{k+1})$ , we can transform the data from the first query model into data from the second query model  $\{W_1, W'_1, \dots, W_{k+1}, W'_{k+1}, \mathbf{w}_1, \mathbf{w}'_1, \dots, \mathbf{w}_{k+1}, \mathbf{w}'_{k+1}\}$ , completing the first part of the induction step.

The induction step for the other direction is similar and omitted.  $\square$

As in the previous work on lower bounds for sequential designs [27, 2, 12], we use a version of Assouad's lemma [30]. Recall that  $\theta_j$  denotes the  $j$ -th component of  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_d)^\top$ .

**Lemma 6.** For any estimator  $\widehat{\boldsymbol{\theta}}$  and any  $\rho > 0$ ,

$$\max_{\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\rho, \rho\}^d} \mathbb{E}_{\boldsymbol{\theta}, \mathcal{V}_k} [\|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|^2] \geq \frac{\rho^2 d}{2} \left(1 - \frac{1}{d} \sum_{j=1}^d \text{TV}(\mathbb{P}_{+j, \boldsymbol{\theta}^*, \mathcal{V}_k}, \mathbb{P}_{-j, \boldsymbol{\theta}^*, \mathcal{V}_k})\right)$$

with

$$\mathbb{P}_{+j, \boldsymbol{\theta}^*, \mathcal{V}_k} := \frac{1}{2^{d-1}} \sum_{\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\rho, \rho\}^d \text{ and } \theta_j = \theta_j^* + \rho} P_{\boldsymbol{\theta}, \mathcal{V}_k},$$

and

$$\mathbb{P}_{-j, \boldsymbol{\theta}^*, \mathcal{V}_k} := \frac{1}{2^{d-1}} \sum_{\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\rho, \rho\}^d \text{ and } \theta_j = \theta_j^* - \rho} P_{\boldsymbol{\theta}, \mathcal{V}_k}.$$

*Proof.* Let  $\xi$  be either  $+1$  or  $-1$ . If  $\boldsymbol{\theta}$  is an element of  $\boldsymbol{\theta}^* + \{-\rho, \rho\}^d$  with  $\theta_j = \theta_j^* + \xi\rho$ , then,  $(\widehat{\theta}_j - \theta_j)^2 \geq \rho^2 \mathbf{1}(\text{sign}(\widehat{\theta}_j - \theta_j^*) \neq \xi)$ , where the sign function evaluated at zero is defined as  $+1$ . For probability measures  $P, Q$  and a measurable set  $A$ , we have  $P(A) + Q(A^c) = 1 + Q(A^c) - P(A^c) \geq 1 - \text{TV}(P, Q)$ . Rewriting  $\|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|^2 = \sum_{j=1}^d (\widehat{\theta}_j - \theta_j)^2$ , and applying the lower bound

$$\max_{\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\rho, \rho\}^d} \geq \frac{1}{2^d} \sum_{\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\rho, \rho\}^d} = \frac{1}{2^d} \sum_{\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\rho, \rho\}^d \text{ and } \theta_j = \theta_j^* + \rho} + \frac{1}{2^d} \sum_{\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\rho, \rho\}^d \text{ and } \theta_j = \theta_j^* - \rho},$$

we find

$$\begin{aligned} & \max_{\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\rho, \rho\}^d} \mathbb{E}_{\boldsymbol{\theta}, \mathcal{V}_k} [\|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|^2] \\ &= \max_{\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\rho, \rho\}^d} \sum_{j=1}^d \mathbb{E}_{\boldsymbol{\theta}, \mathcal{V}_k} [(\widehat{\theta}_j - \theta_j)^2] \\ &\geq \rho^2 \sum_{j=1}^d \frac{1}{2} \mathbb{P}_{+j, \boldsymbol{\theta}^*, \mathcal{V}_k}(\text{sign}(\widehat{\theta}_j - \theta_j^*) \neq 1) + \frac{1}{2} \mathbb{P}_{-j, \boldsymbol{\theta}^*, \mathcal{V}_k}(\text{sign}(\widehat{\theta}_j - \theta_j^*) \neq -1) \\ &\geq \frac{\rho^2 d}{2} \left(1 - \frac{1}{d} \sum_{j=1}^d \text{TV}(\mathbb{P}_{+j, \boldsymbol{\theta}^*, \mathcal{V}_k}, \mathbb{P}_{-j, \boldsymbol{\theta}^*, \mathcal{V}_k})\right). \end{aligned} \tag{36}$$

The right hand side does not depend anymore on the estimator  $\widehat{\boldsymbol{\theta}}$ .  $\square$

*Proof of Theorem 2.* For  $\sigma^2 = 0$ , there is nothing to prove and therefore, we assume  $\sigma^2 > 0$ .

We work in the equivalent model (16). One can assume that  $\mathbf{v}'_\ell \neq 0$  for all  $\ell = 1, \dots, k$ . Otherwise,  $Z'_\ell = 0$  and changing  $\mathbf{v}'_\ell$  to a non-zero vector yields additional information about the model. We can also assume that for any  $\ell = 1, \dots, k$ ,

$$\|\mathbf{v}'_\ell\| = 1, \tag{37}$$

as any other scaling would simply scale the observation  $Z'_\ell$ . Define

$$\tau := \frac{R}{9\sqrt{d}} \left(1 \wedge \frac{d\sigma}{\sqrt{k}R}\right). \tag{38}$$

One can think of  $\tau^2$  as the convergence rate of an individual component  $\mathbb{E}_{\boldsymbol{\theta}, \nu_k} (\hat{\boldsymbol{\theta}}_j - \boldsymbol{\theta}_j)^2$ .

We want to apply Assouad's lower bound in Lemma 6. For any  $\boldsymbol{\theta}^* \in B_{R/2}(0)$  and any  $\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\tau, \tau\}^d$ , we have  $\|\boldsymbol{\theta}\| \leq \|\boldsymbol{\theta}^*\| + R/9 \leq R$ . Thus, we can lower bound  $\sup_{\boldsymbol{\theta} \in B_R(0)} \geq \sup_{\boldsymbol{\theta}^* \in B_{R/2}(0)} \max_{\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\tau, \tau\}^d}$ . Lemma 6 with  $\rho$  replaced by  $\tau$  gives therefore

$$\sup_{\boldsymbol{\theta} \in B_R(0)} \mathbb{E}_{\boldsymbol{\theta}, \nu_k} [\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|^2] \geq \sup_{\boldsymbol{\theta}^* \in B_{R/2}(0)} \frac{\tau^2 d}{2} \left(1 - \frac{1}{d} \sum_{j=1}^d \text{TV}(\mathbb{P}_{+j, \boldsymbol{\theta}^*, \nu_k}, \mathbb{P}_{-j, \boldsymbol{\theta}^*, \nu_k})\right) \quad (39)$$

with

$$\mathbb{P}_{+j, \boldsymbol{\theta}^*, \nu_k} := \frac{1}{2^{d-1}} \sum_{\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\tau, \tau\}^d \text{ and } \theta_j = \theta_j^* + \tau} P_{\boldsymbol{\theta}, \nu_k}, \quad (40)$$

and

$$\mathbb{P}_{-j, \boldsymbol{\theta}^*, \nu_k} := \frac{1}{2^{d-1}} \sum_{\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\tau, \tau\}^d \text{ and } \theta_j = \theta_j^* - \tau} P_{\boldsymbol{\theta}, \nu_k}. \quad (41)$$

Pinsker's inequality  $\text{TV}(P, Q) \leq \sqrt{\text{KL}(P, Q)/2}$  combined with Cauchy-Schwarz inequality  $(\frac{1}{d} \sum_{j=1}^d b_j)^2 \leq (\sum_{j=1}^d d^{-2})(\sum_{j=1}^d b_j^2) = \frac{1}{d} \sum_{j=1}^d b_j^2$  and the joint convexity of the Kullback-Leibler divergence  $\text{KL}(\lambda_1 P_1 + \dots + \lambda_m P_m, \lambda_1 Q_1 + \dots + \lambda_m Q_m) \leq \sum_{s=1}^m \lambda_s \text{KL}(P_s, Q_s)$  for  $\lambda_1, \dots, \lambda_m \geq 0$  and  $\sum_{s=1}^m \lambda_s = 1$ , yield

$$\begin{aligned} & \left( \frac{1}{d} \sum_{j=1}^d \text{TV}(\mathbb{P}_{+j, \boldsymbol{\theta}^*, \nu_k}, \mathbb{P}_{-j, \boldsymbol{\theta}^*, \nu_k}) \right)^2 \\ & \leq \frac{1}{d} \sum_{j=1}^d \text{TV}(\mathbb{P}_{+j, \boldsymbol{\theta}^*, \nu_k}, \mathbb{P}_{-j, \boldsymbol{\theta}^*, \nu_k})^2 \\ & \leq \frac{1}{2d} \sum_{j=1}^d \text{KL}(\mathbb{P}_{+j, \boldsymbol{\theta}^*, \nu_k}, \mathbb{P}_{-j, \boldsymbol{\theta}^*, \nu_k}) \\ & \leq \frac{1}{2d} \sum_{j=1}^d \frac{1}{2^{d-1}} \sum_{\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\tau, \tau\}^d \text{ and } \theta_j = \theta_j^* + \tau} \text{KL}(P_{\boldsymbol{\theta}, \nu_k}, P_{\boldsymbol{\theta} - 2\tau \mathbf{e}_j, \nu_k}) \\ & = \frac{1}{2^d d} \sum_{\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\tau, \tau\}^d} \sum_{j=1}^d \text{KL}(P_{\boldsymbol{\theta}, \nu_k}, P_{\boldsymbol{\theta} - 2\tau \mathbf{e}_j, \nu_k}) \mathbf{1}(\theta_j = \theta_j^* + \tau). \end{aligned} \quad (42)$$

with  $\mathbf{e}_j$  the  $j$ -th standard basis vector.

In a next step, we need to bound  $\text{KL}(P_{\boldsymbol{\theta}, \nu_k}, P_{\boldsymbol{\theta} - 2\tau \mathbf{e}_j, \nu_k})$ . The chain rule for the Kullback-Leibler divergence states that

$$\begin{aligned} \text{KL}(P_{U, V}, Q_{U, V}) &= \mathbb{E}_{P_V} [\text{KL}(P_{U|V}, Q_{U|V})] + \text{KL}(P_V, Q_V) \\ &= \mathbb{E}_{P_{U, V}} [\text{KL}(P_{U|V}, Q_{U|V}) + \text{KL}(P_V, Q_V)]. \end{aligned}$$

The data are generated sequentially,

$$(\mathbf{v}_1, \mathbf{v}'_1) \rightarrow (Z_1, Z'_1) \rightarrow \dots \rightarrow (\mathbf{v}_k, \mathbf{v}'_k) \rightarrow (Z_k, Z'_k).$$

After the query vectors  $(\mathbf{v}_\ell, \mathbf{v}'_\ell)$  are selected, the queries are given by  $Z_\ell = Y_\ell - \mathbf{X}_\ell^\top \mathbf{v}_\ell$  and  $Z'_\ell = \mathbf{X}'_\ell^\top \mathbf{v}'_\ell$ . Therefore,

$$(Z_\ell, Z'_\ell) | (Z_1, Z'_1, \dots, Z_{\ell-1}, Z'_{\ell-1}, \mathbf{v}_1, \mathbf{v}'_1, \dots, \mathbf{v}_\ell, \mathbf{v}'_\ell) = (Z_\ell, Z'_\ell) | (\mathbf{v}_\ell, \mathbf{v}'_\ell).$$

The distribution of  $(\mathbf{v}_\ell, \mathbf{v}'_\ell) | (Z_1, Z'_1, \dots, Z_{\ell-1}, Z'_{\ell-1}, \mathbf{v}_1, \mathbf{v}'_1, \dots, \mathbf{v}_{\ell-1}, \mathbf{v}'_{\ell-1})$  does not depend on the unknown regression vector  $\boldsymbol{\theta}$ . Write  $Q_\ell$  to denote this distribution. If  $P_{\boldsymbol{\theta}, (Z_\ell, Z'_\ell) | (\mathbf{v}_\ell, \mathbf{v}'_\ell)}$  denotes the distribution of  $(Z_\ell, Z'_\ell) | (\mathbf{v}_\ell, \mathbf{v}'_\ell)$  for the data generating parameter  $\boldsymbol{\theta}$ , the chain rule and the arguments above yield

$$\begin{aligned} \text{KL}(P_{\boldsymbol{\theta}, \nu_k}, P_{\boldsymbol{\theta}', \nu_k}) &= \mathbb{E}_{\boldsymbol{\theta}, \nu_k} \left[ \sum_{\ell=1}^k \text{KL}(P_{\boldsymbol{\theta}, (Z_\ell, Z'_\ell) | (\mathbf{v}_\ell, \mathbf{v}'_\ell)}, P_{\boldsymbol{\theta}', (Z_\ell, Z'_\ell) | (\mathbf{v}_\ell, \mathbf{v}'_\ell)}) + \text{KL}(Q_\ell, Q_\ell) \right] \\ &= \mathbb{E}_{\boldsymbol{\theta}, \nu_k} \left[ \sum_{\ell=1}^k \text{KL}(P_{\boldsymbol{\theta}, (Z_\ell, Z'_\ell) | (\mathbf{v}_\ell, \mathbf{v}'_\ell)}, P_{\boldsymbol{\theta}', (Z_\ell, Z'_\ell) | (\mathbf{v}_\ell, \mathbf{v}'_\ell)}) \right]. \end{aligned} \quad (43)$$

We will now apply this identity for  $\boldsymbol{\theta}' = \boldsymbol{\theta} - 2\tau \mathbf{e}_j$  to eventually derive a bound for (42).

For two centered  $d$ -variate normal distributions  $\mathcal{N}(0, \Sigma_0)$ ,  $\mathcal{N}(0, \Sigma_1)$ , with  $\Sigma_0, \Sigma_1 > 0$ , we have

$$\text{KL}(\mathcal{N}(0, \Sigma_1^2), \mathcal{N}(0, \Sigma_0^2)) = \frac{1}{2} \left( \log \left( \frac{\det(\Sigma_0)}{\det(\Sigma_1)} \right) + \text{tr}(\Sigma_1 \Sigma_0^{-1}) - d \right).$$

Consider an invertible matrix  $T$ , and let  $S_0 := T \Sigma_0 T^\top$ ,  $S_1 := T \Sigma_1 T^\top$ . Assume further that there are positive definite matrices  $\Lambda_0, \Lambda_1$  such that  $S_0 \geq \Lambda_0$  and  $S_1 \geq \Lambda_1$ . It is known (e.g. Theorem 6.8 in [31]), that  $S_0^{-1} \leq \Lambda_0^{-1}$ ,  $S_1^{-1} \leq \Lambda_1^{-1}$ , and that  $\text{tr}(AB) = \text{tr}(BA)$  for square matrices  $A, B$ . Using these facts, the symmetrized Kullback-Leibler divergence can be bounded as follows

$$\begin{aligned} \text{KL}(\mathcal{N}(0, \Sigma_0^2), \mathcal{N}(0, \Sigma_1^2)) &\leq \text{KL}(\mathcal{N}(0, \Sigma_0^2), \mathcal{N}(0, \Sigma_1^2)) + \text{KL}(\mathcal{N}(0, \Sigma_1^2), \mathcal{N}(0, \Sigma_0^2)) \\ &= \frac{1}{2} \left( \text{tr}(\Sigma_1 \Sigma_0^{-1}) + \text{tr}(\Sigma_0 \Sigma_1^{-1}) - 2d \right) \\ &= \frac{1}{2} \text{tr}(\Sigma_1^{-1} (\Sigma_1 - \Sigma_0) \Sigma_0^{-1} (\Sigma_1 - \Sigma_0)) \\ &= \frac{1}{2} \text{tr}((T \Sigma_1 T^\top)^{-1} (T \Sigma_1 T^\top - T \Sigma_0 T^\top) (T \Sigma_0 T^\top)^{-1} (T \Sigma_1 T^\top - T \Sigma_0 T^\top)) \\ &= \frac{1}{2} \text{tr}(S_1^{-1} (S_1 - S_0) S_0^{-1} (S_1 - S_0)) \\ &= \frac{1}{2} \text{tr}(S_1^{-1/2} (S_1 - S_0) S_0^{-1} (S_1 - S_0) S_1^{-1/2}) \\ &\leq \frac{1}{2} \text{tr}(S_1^{-1/2} (S_1 - S_0) \Lambda_0^{-1} (S_1 - S_0) S_1^{-1/2}) \\ &= \frac{1}{2} \text{tr}(\Lambda_0^{-1/2} (S_1 - S_0) S_1^{-1} (S_1 - S_0) \Lambda_0^{-1/2}) \\ &\leq \frac{1}{2} \text{tr}(\Lambda_0^{-1/2} (S_1 - S_0) \Lambda_1^{-1} (S_1 - S_0) \Lambda_0^{-1/2}) \\ &= \frac{1}{2} \text{tr}(\Lambda_0^{-1} (S_1 - S_0) \Lambda_1^{-1} (S_1 - S_0)). \end{aligned} \quad (44)$$

We now control the right hand side of (43) for  $\boldsymbol{\theta}' = \boldsymbol{\theta} - 2\tau \mathbf{e}_j$ . As we work in the equivalent model (16) with the normalization constraint (37), the distributions of  $(Z_\ell, Z'_\ell)^\top | (\mathbf{v}_\ell, \mathbf{v}'_\ell)$  is given by (17).



We have normalized  $\mathbf{v}'_\ell$  to a unit-length vector  $\|\mathbf{v}'_\ell\| = 1$ . Thus, the distributions  $P_{\boldsymbol{\theta}, (Z_\ell, Z'_\ell) | (\mathbf{v}_\ell, \mathbf{v}'_\ell)}$  and  $P_{\boldsymbol{\theta} - 2\tau\mathbf{e}_j, (Z_\ell, Z'_\ell) | (\mathbf{v}_\ell, \mathbf{v}'_\ell)}$  (as shown in (17)) are centered normal with respective covariances

$$\begin{pmatrix} \sigma^2 + \|\boldsymbol{\theta} - \mathbf{v}_\ell\|^2 & \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \\ \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle & 1 \end{pmatrix}$$

and

$$\begin{pmatrix} \sigma^2 + \|\boldsymbol{\theta} - 2\tau\mathbf{e}_j - \mathbf{v}_\ell\|^2 & \langle \boldsymbol{\theta} - 2\tau\mathbf{e}_j - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \\ \langle \boldsymbol{\theta} - 2\tau\mathbf{e}_j - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle & 1 \end{pmatrix}.$$

In a next step, we transform the model such that the subsequent analysis becomes more tractable. Choosing the transformation

$$T = \begin{pmatrix} 1 & -\langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \\ 0 & 1 \end{pmatrix}$$

maps  $(Z_\ell, Z'_\ell)^\top$  to  $(Z_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle Z'_\ell, Z'_\ell)^\top$  transforming the covariances via the identities

$$S_0 := T \begin{pmatrix} \sigma^2 + \|\boldsymbol{\theta} - \mathbf{v}_\ell\|^2 & \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \\ \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle & 1 \end{pmatrix} T^\top = \begin{pmatrix} \sigma^2 + \|\boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2 & 0 \\ 0 & 1 \end{pmatrix} \quad (45)$$

and

$$\begin{aligned} S_1 &:= T \begin{pmatrix} \sigma^2 + \|\boldsymbol{\theta} - 2\tau\mathbf{e}_j - \mathbf{v}_\ell\|^2 & \langle \boldsymbol{\theta} - 2\tau\mathbf{e}_j - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \\ \langle \boldsymbol{\theta} - 2\tau\mathbf{e}_j - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle & 1 \end{pmatrix} T^\top \\ &= \begin{pmatrix} \sigma^2 + \|\boldsymbol{\theta} - 2\tau\mathbf{e}_j - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2 & -2\tau \langle \mathbf{e}_j, \mathbf{v}'_\ell \rangle \\ -2\tau \langle \mathbf{e}_j, \mathbf{v}'_\ell \rangle & 1 \end{pmatrix}. \end{aligned} \quad (46)$$

For a symmetric  $2 \times 2$  matrix, the elementary inequality  $2ab \geq -|2ab| \geq -2a^2 - \frac{1}{2}b^2$  yields for all vectors  $(u_1, u_2)^\top$ ,

$$(u_1 u_2)^\top \begin{pmatrix} \alpha & \beta \\ \beta & 1 \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} = u_1^2 \alpha + 2u_1 u_2 \beta + u_2^2 \geq u_1^2 (\alpha - 2\beta^2) + \frac{1}{2} u_2^2$$

implying the matrix inequality

$$\begin{pmatrix} \alpha & \beta \\ \beta & 1 \end{pmatrix} \geq \begin{pmatrix} \alpha - 2\beta^2 & 0 \\ 0 & \frac{1}{2} \end{pmatrix}. \quad (47)$$

If  $k \geq d^2$ , then

$$\tau = \frac{R}{9\sqrt{d}} \left( 1 \wedge \frac{d\sigma}{\sqrt{k}R} \right) \leq \frac{\sigma}{9} \sqrt{\frac{d}{k}} \leq \frac{\sigma}{4}. \quad (48)$$

The last inequality is rather loose but sufficient to give  $8\tau^2 \leq \sigma^2/2$ . Together with the matrix inequality applied to  $\beta = -2\tau \langle \mathbf{e}_j, \mathbf{v}'_\ell \rangle$ , we obtain

$$S_1 \geq \begin{pmatrix} \frac{1}{2}\sigma^2 + \|\boldsymbol{\theta} - 2\tau\mathbf{e}_j - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2 & 0 \\ 0 & \frac{1}{2} \end{pmatrix} \geq \frac{1}{2} \begin{pmatrix} \sigma^2 & 0 \\ 0 & 1 \end{pmatrix} =: \Lambda_1.$$

Setting  $\Lambda_0 := S_0$ , we can now apply (44) for these choices of  $S_0, S_1, \Lambda_0, \Lambda_1$ . Observing that

$$\begin{aligned} & \|\boldsymbol{\theta} - 2\tau\mathbf{e}_j - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2 - \|\boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2 \\ &= 4\tau^2 - 4\tau\langle \mathbf{e}_j, \boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell \rangle, \end{aligned}$$

$$S_1 - S_0 = \begin{pmatrix} 4\tau^2 - 4\tau\langle \mathbf{e}_j, \boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell \rangle & -2\tau\langle \mathbf{e}_j, \mathbf{v}'_\ell \rangle \\ -2\tau\langle \mathbf{e}_j, \mathbf{v}'_\ell \rangle & 0 \end{pmatrix},$$

and

$$\frac{1}{2} \text{tr} \left( \begin{pmatrix} \lambda_0^{-1} & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \alpha & \beta \\ \beta & 0 \end{pmatrix} 2 \begin{pmatrix} \lambda_1^{-1} & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \alpha & \beta \\ \beta & 0 \end{pmatrix} \right) \quad (49)$$

$$= \text{tr} \left( \begin{pmatrix} \alpha\lambda_0^{-1} & \beta\lambda_0^{-1} \\ \beta & 0 \end{pmatrix} \begin{pmatrix} \alpha\lambda_1^{-1} & \beta\lambda_1^{-1} \\ \beta & 0 \end{pmatrix} \right) \quad (50)$$

$$= \text{tr} \left( \begin{pmatrix} \alpha^2\lambda_0^{-1}\lambda_1^{-1} + \beta^2\lambda_0^{-1} & \alpha\beta\lambda_0^{-1}\lambda_1^{-1} \\ \alpha\beta\lambda_1^{-1} & \beta^2\lambda_1^{-1} \end{pmatrix} \right) \quad (51)$$

$$= \frac{\alpha^2}{\lambda_0\lambda_1} + \beta^2 \left( \frac{1}{\lambda_0} + \frac{1}{\lambda_1} \right), \quad (52)$$

we find for the specific choices  $\lambda_0 = \sigma^2 + \|\boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2$ ,  $\lambda_1 = \sigma^2$ ,  $\alpha^2 = (4\tau^2 - 4\tau\langle \mathbf{e}_j, \boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell \rangle)^2 \leq 32\tau^4 + 32\tau^2\langle \mathbf{e}_j, \boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell \rangle^2$  and  $\beta = -2\tau\langle \mathbf{e}_j, \mathbf{v}'_\ell \rangle$ , using (44), and  $\lambda_0 \geq \sigma^2$ ,

$$\begin{aligned} & \text{KL} \left( P_{\boldsymbol{\theta}, (Z_\ell, Z'_\ell) | (\mathbf{v}_\ell, \mathbf{v}'_\ell)}, P_{\boldsymbol{\theta} - 2\tau\mathbf{e}_j, (Z_\ell, Z'_\ell) | (\mathbf{v}_\ell, \mathbf{v}'_\ell)} \right) \\ & \leq \frac{1}{2} \text{tr} \left( \Lambda_0^{-1} (S_1 - S_0) \Lambda_1^{-1} (S_1 - S_0) \right) \\ & = \frac{\alpha^2}{\lambda_0\lambda_1} + \beta^2 \left( \frac{1}{\lambda_0} + \frac{1}{\lambda_1} \right) \\ & \leq \frac{32\tau^4 + 32\tau^2\langle \mathbf{e}_j, \boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell \rangle^2}{(\sigma^2 + \|\boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2)\sigma^2} + 6\frac{\tau^2}{\sigma^2}\langle \mathbf{e}_j, \mathbf{v}'_\ell \rangle^2 \\ & \leq 32\frac{\tau^4}{\sigma^4} + 32\frac{\tau^2\langle \mathbf{e}_j, \boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell \rangle^2}{\sigma^2\|\boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2} + 6\frac{\tau^2}{\sigma^2}\langle \mathbf{e}_j, \mathbf{v}'_\ell \rangle^2. \end{aligned}$$

Interchanging the sums over  $j$  and  $\ell$ , we conclude from (42), (43), the previous inequality,  $\sum_{j=1}^d \langle \mathbf{e}_j, \mathbf{v}'_\ell \rangle^2 = \|\mathbf{v}'_\ell\|^2$ ,  $\sum_{j=1}^d \langle \mathbf{e}_j, \boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell \rangle^2 = \|\boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2$ , the

normalization  $\|\mathbf{v}'_\ell\| = 1$ ,  $\tau \leq \sigma\sqrt{d/(81k)}$  (as derived in (48)), and  $k \geq d^2$ ,

$$\begin{aligned}
& \left( \frac{1}{d} \sum_{j=1}^d \text{TV}(\mathbb{P}_{+j, \boldsymbol{\theta}^*, \mathcal{V}_k}, \mathbb{P}_{-j, \boldsymbol{\theta}^*, \mathcal{V}_k}) \right)^2 \\
& \leq \frac{1}{d} \sum_{j=1}^d \text{TV}(\mathbb{P}_{+j, \boldsymbol{\theta}^*, \mathcal{V}_k}, \mathbb{P}_{-j, \boldsymbol{\theta}^*, \mathcal{V}_k})^2 \\
& \leq \frac{1}{2^d d} \sum_{\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\tau, \tau\}^d} \sum_{j=1}^d \text{KL}(P_{\boldsymbol{\theta}, \mathcal{V}_k}, P_{\boldsymbol{\theta} - 2\tau \mathbf{e}_j, \mathcal{V}_k}) \mathbf{1}(\theta_j = \theta_j^* + \tau) \\
& \leq \frac{1}{2^d d} \sum_{\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\tau, \tau\}^d} \sum_{j=1}^d \sum_{\ell=1}^k 32 \frac{\tau^4}{\sigma^4} + 32 \frac{\tau^2 \langle \mathbf{e}_j, \boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell \rangle^2}{\sigma^2 \|\boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2} + 6 \frac{\tau^2}{\sigma^2} \langle \mathbf{e}_j, \mathbf{v}'_\ell \rangle^2 \\
& = \frac{1}{2^d d} \sum_{\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\tau, \tau\}^d} \sum_{\ell=1}^k 32d \frac{\tau^4}{\sigma^4} + 38 \frac{\tau^2}{\sigma^2} \\
& = 32k \frac{\tau^4}{\sigma^4} + 38 \frac{\tau^2 k}{\sigma^2 d} \\
& \leq \frac{32d^2}{81^2 k} + \frac{38}{81} \\
& \leq \frac{1}{2}.
\end{aligned}$$

Taking square roots, gives  $1 - d^{-1} \sum_{j=1}^d \text{TV}(\mathbb{P}_{+j, \boldsymbol{\theta}^*, \mathcal{V}_k}, \mathbb{P}_{-j, \boldsymbol{\theta}^*, \mathcal{V}_k}) \geq 1 - 1/\sqrt{2}$ . Since by definition  $\tau = \frac{R}{9\sqrt{d}}(1 \wedge \frac{d\sigma}{\sqrt{k}R})$ , the claimed lower bound follows from (39),

$$\begin{aligned}
\sup_{\boldsymbol{\theta} \in \Theta} \mathbb{E}_{\boldsymbol{\theta}, \mathcal{V}_k} [\|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|^2] & \geq \frac{\tau^2 d}{2} \sup_{\boldsymbol{\theta}^* \in B_{R/2}(0)} \left( 1 - \frac{1}{d} \sum_{j=1}^d \text{TV}(\mathbb{P}_{+j, \boldsymbol{\theta}^*, \mathcal{V}_k}, \mathbb{P}_{-j, \boldsymbol{\theta}^*, \mathcal{V}_k}) \right) \\
& = \frac{1}{162} \left( 1 - \frac{1}{\sqrt{2}} \right) \left( R^2 \wedge \frac{d^2}{k} \sigma^2 \right).
\end{aligned}$$

□

*Proof of Theorem 4.* Let

$$\rho = c \frac{R}{\sqrt{d}} \left( 1 \wedge \frac{d}{\sqrt{k}} \left( 1 \vee \frac{\sigma}{R} \right) \right) \quad \text{for } c = 2^{-8}. \quad (53)$$

For any  $\boldsymbol{\theta}^* \in B_{R/2}(0)$  and any  $\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\rho, \rho\}^d$ , we have  $\|\boldsymbol{\theta}\| \leq \|\boldsymbol{\theta}^*\| + cR \leq R$ . Thus, we can lower bound  $\sup_{\boldsymbol{\theta} \in B_R(0)} \geq \sup_{\boldsymbol{\theta}^* \in B_{R/2}(0)} \max_{\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\rho, \rho\}^d}$ . Lemma 6 gives

$$\sup_{\boldsymbol{\theta} \in B_R(0)} \mathbb{E}_{\boldsymbol{\theta}, \mathcal{V}_k} [\|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|^2] \geq \sup_{\boldsymbol{\theta}^* \in B_{R/2}(0)} \frac{\rho^2 d}{2} \left( 1 - \frac{1}{d} \sum_{j=1}^d \text{TV}(\mathbb{P}_{+j, \boldsymbol{\theta}^*, \mathcal{V}_k}, \mathbb{P}_{-j, \boldsymbol{\theta}^*, \mathcal{V}_k}) \right) \quad (54)$$

with  $\mathbb{P}_{+j, \boldsymbol{\theta}^*, \mathcal{V}_k}, \mathbb{P}_{-j, \boldsymbol{\theta}^*, \mathcal{V}_k}$  as in (40) and (41) with  $\tau$  replaced by  $\rho$ .

Since the query vectors are deterministic, (43) becomes now

$$\text{KL}(P_{\boldsymbol{\theta}, \mathcal{V}_k}, P_{\boldsymbol{\theta} - 2\rho \mathbf{e}_j, \mathcal{V}_k}) = \sum_{\ell=1}^k \text{KL}(P_{\boldsymbol{\theta}, (Z_\ell, Z'_\ell) | (\mathbf{v}_\ell, \mathbf{v}'_\ell)}, P_{\boldsymbol{\theta} - 2\rho \mathbf{e}_j, (Z_\ell, Z'_\ell) | (\mathbf{v}_\ell, \mathbf{v}'_\ell)}), \quad (55)$$

where  $\boldsymbol{\theta}$  is an element in the set  $\boldsymbol{\theta}^* + \{-\rho, \rho\}^d$  with  $\theta_j = \theta_j^* + \rho$ .

Let  $\mathbf{v}'$  satisfy  $\|\mathbf{v}'\| = 1$  and let  $\mathbf{w}, \mathbf{w}' \in \{-1, 0, 1\}^d$ . We now prove that  $\|\mathbf{u}\| \geq Rk^{-1/(d-1)}/4$  implies

$$\sigma^2 + \|\mathbf{u} + \rho\mathbf{w} + \rho\langle \mathbf{w}', \mathbf{v}' \rangle \mathbf{v}'\|^2 - 8\rho^2 \geq \frac{\sigma^2 \vee \|\mathbf{u}\|^2}{8}. \quad (56)$$

To see this, use that  $\|\mathbf{v}'\| = 1$  and  $|\langle \mathbf{w}', \mathbf{v}' \rangle| \leq \|\mathbf{w}'\|$ . Thus, triangle inequality yields  $\rho\|\mathbf{w} + \langle \mathbf{w}', \mathbf{v}' \rangle \mathbf{v}'\| \leq 2\rho\sqrt{d}$ . By (53),  $\rho = cRd^{-1/2}(1 \wedge dk^{-1/2}(1 \vee (\sigma/R)))$ .

If  $k \leq 1/(16c)^{d-1}$ , we have  $\|\mathbf{u}\| \geq 4Rc \geq 4\rho\sqrt{d}$ . Moreover, if  $k > 1/(16c)^{d-1}$  and  $\sigma \leq R$ , we can write  $k = (a/16c)^{d-1}$  with  $a > 1$ . Since  $y \leq e^{y-1}$  and  $c \leq 1/(16e^2)$ , we have  $d \leq (e^2)^{(d-1)/2} \leq 1/(16c)^{(d-1)/2}$ . Recalling that  $d \geq 3$ , we find again  $\|\mathbf{u}\| \geq Rk^{-1/(d-1)}/4 = 4cR/a \geq 4cRd/(a^{(d-1)/2}(1/16c)^{(d-1)/2}) = 4cRd/\sqrt{k} \geq 4\rho\sqrt{d}$ . Thus in both of the previous cases  $k \leq 1/(16c)^{d-1}$  and  $\{k > 1/(16c)^{d-1}\} \cap \{\sigma \leq R\}$ ,  $\|\mathbf{u}\| \geq 4\rho\sqrt{d} \geq 2\rho\|\mathbf{w} + \langle \mathbf{w}', \mathbf{v}' \rangle \mathbf{v}'\|$  and therefore by triangle inequality,  $\|\mathbf{u} + \rho\mathbf{w} + \rho\langle \mathbf{w}', \mathbf{v}' \rangle \mathbf{v}'\| \geq \|\mathbf{u}\|/2$ . Since for  $d \geq 4$ , also  $\|\mathbf{u}\|^2/8 \geq 2\rho^2d \geq 8\rho^2$ , (56) follows for these cases.

The remaining case  $k > 1/(16c)^{d-1}$  and  $\sigma > R$ , yields  $\sigma/4 \geq 2cdk^{-1/2}\sigma \geq 2\rho\sqrt{d}$ . Together with the elementary inequality  $a^2 + b^2 \geq \frac{1}{2}(a+b)^2$ , triangle inequality and  $\rho\|\mathbf{w} + \langle \mathbf{w}', \mathbf{v}' \rangle \mathbf{v}'\| \leq 2\rho\sqrt{d}$ , we get  $\sigma + \|\mathbf{u} + \rho\mathbf{w} + \rho\langle \mathbf{w}', \mathbf{v}' \rangle \mathbf{v}'\| \geq \sigma + \|\mathbf{u}\| - 2\rho\sqrt{d} \geq (3/4)\sigma + \|\mathbf{u}\|$  and

$$\sigma^2 + \|\mathbf{u} + \rho\mathbf{w} + \rho\langle \mathbf{w}', \mathbf{v}' \rangle \mathbf{v}'\|^2 \geq \frac{1}{2}(\sigma + \|\mathbf{u} + \rho\mathbf{w} + \rho\langle \mathbf{w}', \mathbf{v}' \rangle \mathbf{v}'\|)^2 \geq \frac{\sigma^2 \vee \|\mathbf{u}\|^2}{4}.$$

Moreover  $\sigma^2/8 \geq 8\rho^2d \geq 8\rho^2$ . Thus, (56) holds in all cases.

Let

$$S_0 = \begin{pmatrix} \sigma^2 + \|\boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2 & 0 \\ 0 & 1 \end{pmatrix}$$

be as in (45).

For the next arguments, we will assume that

$$\|\boldsymbol{\theta}^* - \mathbf{v}_\ell + \langle \boldsymbol{\theta}^* - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle\| \geq \frac{R}{4}k^{-1/(d-1)}. \quad (57)$$

Note that  $\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\rho, \rho\}^d$ . Thus, applying (56) with  $\mathbf{u} = \boldsymbol{\theta}^* - \mathbf{v}_\ell + \langle \boldsymbol{\theta}^* - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle$ ,  $\rho\mathbf{w} = \boldsymbol{\theta} - \boldsymbol{\theta}^*$ , we have that  $S_0 \geq \tilde{\Lambda}_0$  with

$$\tilde{\Lambda}_0 = \begin{pmatrix} \frac{\sigma^2 \vee \|\boldsymbol{\theta}^* - \mathbf{v}_\ell - \langle \boldsymbol{\theta}^* - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2}{8} & 0 \\ 0 & 1 \end{pmatrix}.$$

Similarly, let

$$S_1 = \begin{pmatrix} \sigma^2 + \|\boldsymbol{\theta} - 2\rho\mathbf{e}_j - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2 & -2\rho\langle \mathbf{e}_j, \mathbf{v}'_\ell \rangle \\ -2\rho\langle \mathbf{e}_j, \mathbf{v}'_\ell \rangle & 1 \end{pmatrix}$$

be as in (46). Using the matrix inequality (47),  $\langle \mathbf{e}_j, \mathbf{v}'_\ell \rangle \leq \|\mathbf{e}_j\| \|\mathbf{v}'_\ell\| = 1$ , and applying (56) with  $\mathbf{w} = \boldsymbol{\theta} - \boldsymbol{\theta}^* - 2\rho\mathbf{e}_j$  and  $\mathbf{w}' = \boldsymbol{\theta} - \boldsymbol{\theta}^*$ , we have

$$\begin{aligned} S_1 &\geq \begin{pmatrix} \sigma^2 + \|\boldsymbol{\theta} - 2\rho\mathbf{e}_j - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2 - 8\rho^2 & 0 \\ 0 & \frac{1}{2} \end{pmatrix} \\ &\geq \begin{pmatrix} \frac{\sigma^2 \vee \|\boldsymbol{\theta}^* - \mathbf{v}_\ell - \langle \boldsymbol{\theta}^* - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2}{8} & 0 \\ 0 & \frac{1}{2} \end{pmatrix}. \end{aligned}$$

Recall that  $\boldsymbol{\theta}$  is an element in the set  $\boldsymbol{\theta}^* + \{-\rho, \rho\}^d$  with  $\theta_j = \theta_j^* + \rho$ . This means that  $\langle \boldsymbol{\theta} - \rho\mathbf{e}_j - \boldsymbol{\theta}^*, \mathbf{e}_j \rangle = 0$  and thus

$$\begin{aligned} \|\boldsymbol{\theta} - 2\rho\mathbf{e}_j - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2 &- \|\boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2 \\ &= -4\rho \langle \mathbf{e}_j, \boldsymbol{\theta} - \rho\mathbf{e}_j - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell \rangle \\ &= -4\rho \langle \mathbf{e}_j, \boldsymbol{\theta}^* - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell \rangle. \end{aligned}$$

This implies

$$S_1 - S_0 = \begin{pmatrix} -4\rho \langle \mathbf{e}_j, \boldsymbol{\theta}^* - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell \rangle & -2\rho \langle \mathbf{e}_j, \mathbf{v}'_\ell \rangle \\ -2\rho \langle \mathbf{e}_j, \mathbf{v}'_\ell \rangle & 0 \end{pmatrix} = \begin{pmatrix} \alpha & \beta \\ \beta & 0 \end{pmatrix}.$$

Combining (44) and (52) with  $\lambda_0 = \lambda_1/2 = \frac{1}{8}(\sigma^2 \vee \|\boldsymbol{\theta}^* - \mathbf{v}_\ell - \langle \boldsymbol{\theta}^* - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2)$ , and  $\alpha, \beta$  as above, we obtain

$$\begin{aligned} \text{KL}(P_{\boldsymbol{\theta}, (Z_\ell, Z'_\ell) | (\mathbf{v}_\ell, \mathbf{v}'_\ell)}, P_{\boldsymbol{\theta} - 2\rho\mathbf{e}_j, (Z_\ell, Z'_\ell) | (\mathbf{v}_\ell, \mathbf{v}'_\ell)}) &\leq 32 \frac{(-4\rho \langle \mathbf{e}_j, \boldsymbol{\theta}^* - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell \rangle)^2}{(\sigma^2 \vee \|\boldsymbol{\theta}^* - \mathbf{v}_\ell - \langle \boldsymbol{\theta}^* - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2)^2} \\ &\quad + 12 \frac{(2\rho \langle \mathbf{e}_j, \mathbf{v}'_\ell \rangle)^2}{\sigma^2 \vee \|\boldsymbol{\theta}^* - \mathbf{v}_\ell - \langle \boldsymbol{\theta}^* - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2}. \end{aligned}$$

Since  $\langle \boldsymbol{\theta} - \boldsymbol{\theta}^*, \mathbf{v}'_\ell \rangle^2 \leq \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|^2 \|\mathbf{v}'_\ell\|^2 = \rho^2 d$ , we find using  $(a+b)^2 \leq 2a^2 + 2b^2$ ,

$$\langle \mathbf{e}_j, \boldsymbol{\theta}^* - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell \rangle^2 \leq 2 \langle \mathbf{e}_j, \boldsymbol{\theta}^* - \langle \boldsymbol{\theta}^* - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell \rangle^2 + 2\rho^2 d \langle \mathbf{e}_j, \mathbf{v}'_\ell \rangle^2.$$

Since  $\mathbf{e}_1, \dots, \mathbf{e}_d$  form an orthonormal basis,  $\sum_{j=1}^d \langle \mathbf{e}_j, \mathbf{a} \rangle^2 = \|\mathbf{a}\|^2$  for all  $d$ -dimensional vectors  $\mathbf{a}$ . Thus taking the sum over  $j$  yields

$$\begin{aligned} &\sum_{j=1}^d \text{KL}(P_{\boldsymbol{\theta}, (Z_\ell, Z'_\ell) | (\mathbf{v}_\ell, \mathbf{v}'_\ell)}, P_{\boldsymbol{\theta} - 2\rho\mathbf{e}_j, (Z_\ell, Z'_\ell) | (\mathbf{v}_\ell, \mathbf{v}'_\ell)}) \\ &\leq (2^8 + 48) \frac{\rho^2}{\sigma^2 \vee \|\boldsymbol{\theta}^* - \mathbf{v}_\ell - \langle \boldsymbol{\theta}^* - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2} + 2^8 \frac{\rho^4 d}{(\sigma^2 \vee \|\boldsymbol{\theta}^* - \mathbf{v}_\ell - \langle \boldsymbol{\theta}^* - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2)^2}. \end{aligned}$$

Recall that by (55),  $\text{KL}(P_{\boldsymbol{\theta}, \mathcal{V}_k}, P_{\boldsymbol{\theta} - 2\rho\mathbf{e}_j, \mathcal{V}_k}) = \sum_{\ell=1}^k \text{KL}(P_{\boldsymbol{\theta}, (Z_\ell, Z'_\ell) | (\mathbf{v}_\ell, \mathbf{v}'_\ell)}, P_{\boldsymbol{\theta} - 2\rho\mathbf{e}_j, (Z_\ell, Z'_\ell) | (\mathbf{v}_\ell, \mathbf{v}'_\ell)})$ . Hence,

$$\begin{aligned} &\sum_{j=1}^d \text{KL}(P_{\boldsymbol{\theta}, \mathcal{V}_k}, P_{\boldsymbol{\theta} - 2\rho\mathbf{e}_j, \mathcal{V}_k}) \\ &\leq \sum_{\ell=1}^k \frac{176\rho^2}{\sigma^2 \vee \|\boldsymbol{\theta}^* - \mathbf{v}_\ell - \langle \boldsymbol{\theta}^* - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2} + \sum_{\ell=1}^k \frac{2^8 \rho^4 d}{(\sigma^2 \vee \|\boldsymbol{\theta}^* - \mathbf{v}_\ell - \langle \boldsymbol{\theta}^* - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2)^2}. \end{aligned}$$

The right hand side does not depend on  $\boldsymbol{\theta}$  anymore. Combined with (42),

$$\begin{aligned}
& \left( \frac{1}{d} \sum_{j=1}^d \text{TV} \left( \mathbb{P}_{+j, \boldsymbol{\theta}^*, \mathcal{V}_k}, \mathbb{P}_{-j, \boldsymbol{\theta}^*, \mathcal{V}_k} \right) \right)^2 \\
& \leq \frac{1}{2^d d} \sum_{\boldsymbol{\theta} \in \boldsymbol{\theta}^* + \{-\tau, \tau\}^d} \sum_{j=1}^d \text{KL} \left( P_{\boldsymbol{\theta}, \mathcal{V}_k}, P_{\boldsymbol{\theta} - 2\tau \mathbf{e}_j, \mathcal{V}_k} \right) \mathbf{1}(\theta_j = \theta_j^* + \tau) \\
& \leq \sum_{\ell=1}^k \frac{88\rho^2}{d(\sigma^2 \vee \|\boldsymbol{\theta}^* - \mathbf{v}_\ell - \langle \boldsymbol{\theta}^* - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2)} + \sum_{\ell=1}^k \frac{2^7 \rho^4}{(\sigma^2 \vee \|\boldsymbol{\theta}^* - \mathbf{v}_\ell - \langle \boldsymbol{\theta}^* - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2)^2}.
\end{aligned} \tag{58}$$

Notice that this bound assumes the inequality (57).

If  $\|\mathbf{v}'\| = 1$ , then  $\|\mathbf{u} - \mathbf{v} - \langle \mathbf{u} - \mathbf{v}, \mathbf{v}' \rangle \mathbf{v}'\| = \inf_{\gamma \in \mathbb{R}} \|\mathbf{u} - \mathbf{v} - \gamma \mathbf{v}'\|$ . Thus, the set  $T_{\mathbf{v}, \mathbf{v}'}(r) := \{\mathbf{u} : \|\mathbf{u} - \mathbf{v} - \langle \mathbf{u} - \mathbf{v}, \mathbf{v}' \rangle \mathbf{v}'\| \leq r\}$  is a tube with radius  $r > 0$ . The Euclidean ball with radius  $r' > 0$  and center  $\mathbf{w}$  is denoted by  $B_{r'}(\mathbf{w}) := \{\mathbf{y} : \|\mathbf{y} - \mathbf{w}\| \leq r'\}$ . If  $\mathbf{u} \in B_{r'}(0)$ , then  $\langle -\mathbf{v}, \mathbf{v}' \rangle - r' \leq \langle \mathbf{u} - \mathbf{v}, \mathbf{v}' \rangle \leq \langle -\mathbf{v}, \mathbf{v}' \rangle + r'$ . Thus,

$$H_{\mathbf{v}, \mathbf{v}'}(r, r') := T_{\mathbf{v}, \mathbf{v}'}(r) \cap B_{r'}(0) \subseteq \left\{ \mathbf{u} : \inf_{\langle -\mathbf{v}, \mathbf{v}' \rangle - r' \leq \gamma \leq \langle -\mathbf{v}, \mathbf{v}' \rangle + r'} \|\mathbf{u} - \mathbf{v} - \gamma \mathbf{v}'\| \leq r \right\}.$$

Let  $\Gamma(\cdot)$  be the Gamma function. The volume formula for a tube (e.g. (2.3) in [17]) says that the volume of the right hand side is bounded by  $2r'c_{d-1}r^{d-1} + c_d r^d$  with  $c_p := \pi^{p/2}/\Gamma(p/2 + 1)$  the volume of a  $p$ -dimensional unit ball. By (41) in [23],  $\Gamma(d/2 + 1)/\Gamma(d/2 + 1/2) \leq \sqrt{d/2 + 1/2}$ . Therefore  $c_{d-1} \leq \sqrt{\pi d} c_d$ . Moreover,  $d \geq 5$  implies  $1 + 2\sqrt{\pi} \sqrt{d/2 + 1/2} \leq 2^{d-2}$ . Thus,  $2c_{d-1} + c_d \leq (1 + 2\sqrt{\pi} \sqrt{d/2 + 1/2})c_d \leq 2^{d-2}c_d$ . Together with the previous display, this means that

$$\text{Vol} \left( H_{\mathbf{v}, \mathbf{v}'}(r, r') \right) \leq 2^{d-2} r' c_d r^{d-1}, \quad \text{whenever } r \leq r'.$$

with  $\text{Vol}$  the volume (Lebesgue measure) in  $d$ -dimensions. Using again that  $d \geq 4$ , this means that for any  $\mathbf{v}_1, \mathbf{v}'_1, \dots, \mathbf{v}_k, \mathbf{v}'_k$ , and any  $b > 0$ ,

$$\begin{aligned}
\text{Vol} \left( H_{\mathbf{v}_\ell, \mathbf{v}'_\ell} \left( \frac{bR}{4} k^{-1/(d-1)}, \frac{R}{2} \right) \right) & \leq 2^{d-3} R c_d b^{d-1} R^{d-1} k^{-1} 4^{1-d} \\
& = \frac{b^{d-1}}{2k} \left( \frac{R}{2} \right)^d c_d \\
& = \frac{b^{d-1}}{2k} \text{Vol} \left( B_{R/2}(0) \right)
\end{aligned}$$

and

$$\begin{aligned}
\text{Vol} \left( \bigcup_{\ell=1}^k H_{\mathbf{v}_\ell, \mathbf{v}'_\ell} \left( \frac{R}{4} k^{-1/(d-1)}, \frac{R}{2} \right) \right) & \leq \sum_{\ell=1}^k \text{Vol} \left( H_{\mathbf{v}_\ell, \mathbf{v}'_\ell} \left( \frac{R}{4} k^{-1/(d-1)}, \frac{R}{2} \right) \right) \\
& \leq \frac{1}{2} \text{Vol} \left( B_{R/2}(0) \right).
\end{aligned}$$

The latter implies that

$$\begin{aligned}
\text{Vol} \left( B_{R/2}(0) \setminus \bigcup_{\ell=1}^k H_{\mathbf{v}_\ell, \mathbf{v}'_\ell} \left( \frac{R}{8} k^{-1/(d-1)}, \frac{R}{2} \right) \right) & \geq \text{Vol} \left( B_{R/2}(0) \right) - \frac{1}{2} \text{Vol} \left( B_{R/2}(0) \right) \\
& = \frac{1}{2} \text{Vol} \left( B_{R/2}(0) \right).
\end{aligned}$$

Now we define a probability measure on the ball  $B_{R/2}(0)$  by

$$\nu(A) = \frac{\text{Vol}(A \cap (B_{R/2}(0) \setminus \cup_{\ell=1}^k H_{\mathbf{v}_\ell, \mathbf{v}'_\ell}(Rk^{-1/(d-1)}/4, R/2)))}{\text{Vol}(B_{R/2}(0) \setminus \cup_{\ell=1}^k H_{\mathbf{v}_\ell, \mathbf{v}'_\ell}(Rk^{-1/(d-1)}/4, R/2))}.$$

The distribution ensures that the inequality (57) holds for all  $\ell = 1, \dots, k$  with probability one. Thus, we can apply (58) and therefore lower bound  $\inf_{\boldsymbol{\theta}^* \in B_{R/2}(0)}$  by an average with respect to the probability measure  $\nu$ ,

$$\begin{aligned} & \inf_{\boldsymbol{\theta}^* \in B_{R/2}(0)} \left( \frac{1}{d} \sum_{j=1}^d \text{TV}(\mathbb{P}_{+j, \boldsymbol{\theta}^*, \nu_k}, \mathbb{P}_{-j, \boldsymbol{\theta}^*, \nu_k}) \right)^2 \\ & \leq \int \left( \sum_{\ell=1}^k \frac{88\rho^2}{d(\sigma^2 \vee \|\boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2)} + \sum_{\ell=1}^k \frac{2^7 \rho^4}{(\sigma^2 \vee \|\boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2)^2} \right) d\nu(\boldsymbol{\theta}) \\ & \leq \sum_{\ell=1}^k \frac{88\rho^2}{d} \left( \frac{1}{\sigma^2} \wedge \int \frac{1}{\|\boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2} \right) d\nu(\boldsymbol{\theta}) \\ & \quad + 2^7 \rho^4 \left( \frac{1}{\sigma^4} \wedge \int \frac{1}{\|\boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^4} \right) d\nu(\boldsymbol{\theta}). \end{aligned} \quad (59)$$

The sum  $\sum_{0 \leq t \leq \log_2(k^{1/(d-1)} - 1)}$  is set to zero if  $k < 2^{d-1}$ . Using that for  $d \geq 6$ , one has  $2^{d-5} - 1 \geq \frac{1}{2}2^{d-5}$ , we obtain  $\sum_{t=0}^q 2^{t(d-5)} = (2^{(q+1)(d-5)} - 1)/(2^{d-5} - 1) \leq 2 \cdot 2^{q(d-5)}$ . Applying this with  $q = \lfloor \log_2(k^{1/(d-1)} - 1) \rfloor$  and using that  $k^{-4/(d-1)}k = k^{(d-5)/(d-1)}$  yields

$$\begin{aligned} & \int \frac{1}{\|\boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^4} d\nu(\boldsymbol{\theta}) \\ & \leq \int_{\|\boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2 \geq R/8} \frac{1}{\|\boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^4} d\nu(\boldsymbol{\theta}) \\ & \quad + \frac{2}{\text{Vol}(B_{R/2}(0))} \sum_{0 \leq t \leq \log_2(k^{1/(d-1)} - 1)} \\ & \quad \int_{2^t Rk^{-1/(d-1)}/4 \leq \|\boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2 \leq 2^{t+1} Rk^{-1/(d-1)}/4} \frac{1}{\|\boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^4} d\nu(\boldsymbol{\theta}) \\ & \leq \frac{64}{R^4} + \frac{2}{\text{Vol}(B_{R/2}(0))} \sum_{0 \leq t \leq \log_2(k^{1/(d-1)} - 1)} \frac{\text{Vol}(H_{\mathbf{v}_\ell, \mathbf{v}'_\ell}(2^{t+1} Rk^{-1/(d-1)}/4, R/2))}{(2^{2t} R^2 k^{-2/(d-1)}/16)^2} \\ & = \frac{64}{R^4} + \frac{2}{\text{Vol}(B_{R/2}(0))} \sum_{0 \leq t \leq \log_2(k^{1/(d-1)} - 1)} \frac{2^8}{2^{4t} R^4 k^{-4/(d-1)}} \frac{2^{(t+1)(d-1)}}{2^k} \text{Vol}(B_{R/2}(0)) \\ & \leq \frac{64}{R^4} + \frac{2^8 \cdot 2^5}{R^4} \\ & \leq \frac{2^{14}}{R^4}. \end{aligned} \quad (60)$$

Since  $\nu$  is a probability measure, Jensen's inequality gives

$$\int \frac{1}{\|\boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^2} d\nu(\boldsymbol{\theta}) \leq \left( \int \frac{1}{\|\boldsymbol{\theta} - \mathbf{v}_\ell - \langle \boldsymbol{\theta} - \mathbf{v}_\ell, \mathbf{v}'_\ell \rangle \mathbf{v}'_\ell\|^4} d\nu(\boldsymbol{\theta}) \right)^{1/2} \leq \frac{2^7}{R^2}.$$

Recall that  $\rho = cRd^{-1/2}(1 \wedge dk^{-1/2}(1 \vee \sigma/R))$ . Thus,  $\rho^2 \leq c^2 dk^{-1}(R \vee \sigma)^2$ ,  $\rho^2 \leq c^2 R^2 d^{-1}$ , and

$\rho^4 \leq \rho^2 \rho^2 \leq c^4 k^{-1} (R \vee \sigma)^2 R^2$ . Combined with (59), (60), and  $c = 2^{-8}$ ,

$$\inf_{\boldsymbol{\theta}^* \in B_{R/2}(0)} \left( \frac{1}{d} \sum_{j=1}^d \text{TV}(\mathbb{P}_{+j, \boldsymbol{\theta}^*, \nu_k}, \mathbb{P}_{-j, \boldsymbol{\theta}^*, \nu_k}) \right)^2 \leq \frac{88 \cdot 2^7 k \rho^2}{d(\sigma^2 \vee R^2)} + \frac{2^{21} k \rho^4}{\sigma^4 \vee R^4} \leq 88 \cdot 2^7 c^2 + 2^{21} c^4 \leq \frac{1}{4}.$$

Using (54), we find that

$$\sup_{\boldsymbol{\theta} \in B_R(0)} \mathbb{E}_{\boldsymbol{\theta}, \nu_k} [\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|^2] \geq \frac{\rho^2 d}{4} = 2^{-18} R^2 \left( 1 \wedge \frac{d^2}{k} (R \vee \sigma)^2 \right).$$

□

*Proof of Corollary 1.* Using  $\boldsymbol{\theta}_0 = 0$  and Lemma 4 with  $\mathbf{a} = \mathbf{b} = \boldsymbol{\theta}^*$ , we obtain that  $\|S_0\| = \|\boldsymbol{\theta}^*(\boldsymbol{\theta}^*)^\top\| = \|\boldsymbol{\theta}^*\|^2 \leq d$ . Thus Theorem 1 combined with (14) for  $\kappa = 1$  and  $\gamma = 2$  yields the rate  $d^2 \log^2(d)/k$  for the upper bound in the adaptive query setting.

The minimax rate in the non-adaptive query setting is  $R^2 \wedge \frac{d^2}{k} (R \vee \sigma)^2$ . Using that  $R = \sqrt{d} \geq \sigma$ , this becomes  $d^2 \wedge d^3/k$ . Since  $k \gtrsim d^2 \log(d)$  the second term always dominates and the minimax rate is  $d^3/k$ . □

## Acknowledgement

We are grateful to Chao Gao for fruitful discussions. The research has been supported by the NWO Vidi grant VI.Vidi.192.021.

## References

- [1] AGARWAL, A., DEKEL, O., AND XIAO, L. Optimal algorithms for online convex optimization with multi-point bandit feedback. In *Proceedings of the 23rd Annual Conference on Learning Theory (COLT)* (2010).
- [2] ARIAS-CASTRO, E., CANDES, E. J., AND DAVENPORT, M. A. On the fundamental limits of adaptive sensing. *IEEE Transactions on Information Theory* 59, 1 (2013), 472–481.
- [3] BACH, F., AND MOULINES, E. Non-strongly-convex smooth stochastic approximation with convergence rate  $O(1/n)$ . *Advances in neural information processing systems* 26 (2013).
- [4] BACH, F., AND PERCHET, V. Highly-smooth zero-th order online optimization. In *Conference on Learning Theory* (2016), PMLR, pp. 257–283.
- [5] BALASUBRAMANIAN, K., AND GHADIMI, S. Zeroth-order nonconvex stochastic optimization: handling constraints, high dimensionality, and saddle points. *Found. Comput. Math.* 22, 1 (2022), 35–76.
- [6] BAYDIN, A. G., PEARLMUTTER, B. A., SYME, D., WOOD, F., AND TORR, P. Gradients without backpropagation. *arXiv preprint arXiv:2202.08587* (2022).



- [7] BELLO, K., AND HONORIO, J. Computationally and statistically efficient learning of causal Bayes nets using path queries. In *Advances in Neural Information Processing Systems* (2018), S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., vol. 31, Curran Associates, Inc.
- [8] BOS, T., AND SCHMIDT-HIEBER, J. Convergence guarantees for forward gradient descent in the linear regression model. *arXiv preprint* (2023).
- [9] CLARA, G., LANGER, S., AND SCHMIDT-HIEBER, J. Dropout regularization versus  $\ell_2$ -penalization in the linear model. *arXiv preprint* (2023), arXiv:2306.10529.
- [10] CRICK, F. The recent excitement about neural networks. *Nature* 337, 6203 (1989), 129–132.
- [11] DIEULEVEUT, A., FLAMMARION, N., AND BACH, F. Harder, better, faster, stronger convergence rates for least-squares regression. *The Journal of Machine Learning Research* 18, 1 (2017), 3520–3570.
- [12] DUCHI, J. C., JORDAN, M. I., WAINWRIGHT, M. J., AND WIBISONO, A. Optimal rates for zero-order convex optimization: The power of two function evaluations. *IEEE Transactions on Information Theory* 61, 5 (2015), 2788–2806.
- [13] FLAXMAN, A. D., KALAI, A. T., AND MCMAHAN, H. B. Online convex optimization in the bandit setting: gradient descent without a gradient. *arXiv preprint cs/0408007* (2004).
- [14] FRÉMAUX, N., SPREKELER, H., AND GERSTNER, W. Functional requirements for reward-modulated spike-timing-dependent plasticity. *Journal of Neuroscience* 30, 40 (2010), 13326–13337.
- [15] HAZAN, E., KOREN, T., AND LEVY, K. Y. Logistic regression: Tight bounds for stochastic and online optimization. In *Conference on Learning Theory* (2014), PMLR, pp. 197–209.
- [16] JIN, Y., XIAO, T., AND BALASUBRAMANIAN, K. Statistical inference for Polyak-Ruppert averaged zeroth-order stochastic gradient algorithm. *arXiv preprint* (2021), arXiv:2102.05198.
- [17] JOHNSTONE, I., AND SIEGMUND, D. On Hotelling’s formula for the volume of tubes and Naiman’s inequality. *Ann. Statist.* 17, 1 (1989), 184–194.
- [18] LAKSHMINARAYANAN, C., AND SZEPESVARI, C. Linear stochastic approximation: How far does constant step-size and iterate averaging go? In *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics* (09–11 Apr 2018), A. Storkey and F. Perez-Cruz, Eds., vol. 84 of *Proceedings of Machine Learning Research*, PMLR, pp. 1347–1355.
- [19] LARSON, J., MENICKELLY, M., AND WILD, S. M. Derivative-free optimization methods. *Acta Numer.* 28 (2019), 287–404.
- [20] LILLICRAP, T. P., SANTORO, A., MARRIS, L., AKERMAN, C. J., AND HINTON, G. Backpropagation and the brain. *Nature Reviews Neuroscience* 21, 6 (2020), 335–346.

- [21] LUGOSI, G., TRUSZKOWSKI, J., VELONA, V., AND ZWIERNIK, P. Learning partial correlation graphs and graphical models by covariance queries. *Journal of Machine Learning Research* 22, 203 (2021), 1–41.
- [22] NOVITSKII, V., AND GASNIKOV, A. Improved exploitation of higher order smoothness in derivative-free optimization. *Optimization Letters* 16, 7 (2022), 2059–2071.
- [23] QI, F., AND LUO, Q.-M. Bounds for the ratio of two gamma functions: from Wendel’s asymptotic relation to Elezović-Giordano-Pečarić’s theorem. *J. Inequal. Appl.* (2013), 2013:542, 20.
- [24] REGIS, M., SERRA, P., AND VAN DEN HEUVEL, E. R. Random autoregressive models: a structured overview. *Econometric Reviews* 41, 2 (2022), 207–230.
- [25] REN, M., KORNBLITH, S., LIAO, R., AND HINTON, G. Scaling forward gradient with local losses. *arXiv preprint arXiv:2210.03310* (2022).
- [26] SCHMIDT-HIEBER, J. Interpreting learning in biological neural networks as zero-order optimization method. *arXiv preprint* (2023), arXiv:2301.11777.
- [27] SHAMIR, O. On the complexity of bandit and derivative-free stochastic convex optimization. In *Proceedings of the 26th Annual Conference on Learning Theory* (Princeton, NJ, USA, 2013), S. Shalev-Shwartz and I. Steinwart, Eds., vol. 30 of *Proceedings of Machine Learning Research*, PMLR, pp. 3–24.
- [28] STEINHARDT, J., VALIANT, G., AND WAGER, S. Memory, communication, and statistical queries. In *29th Annual Conference on Learning Theory* (Columbia University, New York, New York, USA, 23–26 Jun 2016), V. Feldman, A. Rakhlin, and O. Shamir, Eds., vol. 49 of *Proceedings of Machine Learning Research*, PMLR, pp. 1490–1516.
- [29] TRAPPENBERG, T. P. *Fundamentals of Computational Neuroscience: Third Edition*. Oxford University Press, 12 2022.
- [30] TSYBAKOV, A. B. *Introduction to nonparametric estimation*. Springer Series in Statistics. Springer, New York, 2009. Revised and extended from the 2004 French original, Translated by Vladimir Zaiats.
- [31] ZHANG, F. *Matrix theory*. Universitext. Springer-Verlag, New York, 1999. Basic results and techniques.