

Posterior concentration and fast convergence rates for generalized Bayesian learning

Lam Si Tung Ho^a, Binh T. Nguyen^{b,d,c}, Vu Dinh^e, Duy Nguyen^f

^a*Department of Mathematics and Statistics, Dalhousie University, Halifax, Nova Scotia, Canada*

^b*Department of Computer Science, Faculty of Mathematics and Computer Science, University of Science, Ho Chi Minh City, Vietnam*

^c*Vietnam National University, Ho Chi Minh City, Vietnam*

^d*AISIA Research Lab*

^e*Department of Mathematical Sciences, University of Delaware, USA*

^f*Department of Statistics, University of Wisconsin-Madison, USA*

Abstract

In this paper, we study the learning rate of generalized Bayes estimators in a general setting where the hypothesis class can be uncountable and have an irregular shape, the loss function can have heavy tails, and the optimal hypothesis may not be unique. We prove that under the multi-scale Bernstein's condition, the generalized posterior distribution concentrates around the set of optimal hypotheses and the generalized Bayes estimator can achieve fast learning rate. Our results are applied to show that the standard Bayesian linear regression is robust to heavy-tailed distributions.

Keywords: Bayesian learning, posterior concentration, fast rate, heavy-tailed loss, Bernstein condition

1. Introduction

There has been a growing interest in posterior concentration rates of Bayesian inference over the last decade. Posterior concentration allows us to uncover frequentist properties of Bayesian methods and implies that most of the posterior mass will be close to the truth in the frequentist sense. Studying such properties enables designs of appropriate priors for Bayesian inference in various contexts [1, 2, 3].

Similar approaches have also been proposed in statistical learning theory. In such settings, one considers models of predictors defined relative to some loss functions and proves frequentist convergence bounds of generalized Bayes predictors constructed with respect to a *posterior randomization measure*. The most notable work on this direction is the framework of “safe Bayesian,” where the formulation for generalized Bayesian posterior can be tuned by an optimal learning rate [4]. Instead of choosing priors, within such a framework, one can construct more flexible estimators over a wide range of hypothesis spaces, losses, and model misspecifications.

From another perspective, the topic of fast learning rate in statistical learning has become a subject of growing interest in recent works. The pursuit of a “fast rate” regime has led to many conditions in learning theory under which fast rates are possible such as low noise assumption [5, 6], stochastic mixability condition [7], Bernstein’s condition [8], v -central condition [9], and multi-scale Bernstein’s condition [10]. Traditionally, most works in this direction have primarily focused on bounded losses, and deviations from this expected behavior are worrisome, especially when the loss of the learning problem of interest is unbounded and/or has heavy tails.

Recently, it has been shown that it is possible to generalize conditions for fast learning rates with unbounded and heavy-tailed losses. The fast learning rate for sub-gaussian and sub-exponential losses are done in the context of density estimation [11, 12] and for general losses [13], of which proofs of fast rates heavily employ the Bernstein’s condition and the central condition. In [14], the authors provide an exponential concentration of the median-of-means estimator under heavy-tailed distributions to approximate minimization of smooth and strongly convex losses. Similarly, the paper [15] proposes studying the “optimistic rate” under the small-ball condition for learning with heavy-tailed convex losses. Another effort to resolve this issue was shown in [10] with their newly proposed multi-scale Bernstein’s condition, which enables learning with heavy tails when the loss function is non-convex and the optimal hypothesis is not unique. Their analyses recover fast learning rates for empirical risk minimization

(ERM) estimators under bounded losses, but, more significantly, also hold for heavy-tailed losses.

The vast majority of the recent works in obtaining fast learning rates have taken place in the frequentist approach, whereas applications to generalized Bayesian estimators are unknown. In [16], the authors take a further step to show that fast learnings in the generalized Bayesian setting are, indeed, attainable. However, the major drawbacks are that the optimal learning rate β must be known in advance and that the hypothesis class is finite. The “safe Bayesian” methods [4] provide a framework to analyze a special form of generalized Bayes estimators employing the central condition, which cannot be applied to losses with polynomial tails [9]. As a result, the feasibility of fast learning rates for heavy-tailed distributions under Bayesian frameworks remains unknown.

Building upon the multi-scale Bernstein’s condition, we analyze fast concentration rates of generalized Bayes estimators in a general framework where the hypothesis class can be infinite/uncountable and have an irregular shape, the loss function can have heavy tails, and the optimal hypothesis may not be unique. Our results demonstrate that learning rates faster than $\mathcal{O}(n^{-1/2})$ can be obtained. Moreover, depending on the regularity of the risk function and the complexity of the hypothesis class, the learning rate can be arbitrarily close to the optimal rate $\mathcal{O}(n^{-1})$. We apply our results to show that the standard Bayesian linear regression is robust to heavy-tailed distributions. Specifically, Bayesian linear regression with the regular square loss can achieve fast rate learning when the errors follow t-distributions.

Related work. Bayesian framework has been applied extensively to a wide variety of research areas including ecology [17], evolutionary biology [18], epidemiology [19, 20], and economics [21]. However, theoretical properties of Bayesian methods have not been studied extensively as its frequentist counterparts, especially for heavy-tailed losses. In particular, several frequentist approaches have been shown to perform well with heavy-tailed losses including ERM [10, 22], median-of-means estimator [14, 23], k-mean clustering [10, 24], support vector machines [25], Least

Squares Estimator [26]. Recently, much effort have been devoted to study the asymptotic theory of Bayesian inference [1, 4, 27, 28, 29, 30]. However, the lack of results for heavy-tailed losses has hindered the applicability of the Bayesian inference to such a scenario. This is a major disadvantage compared to other frequentist methods. Therefore, it is essential to establish a theoretical guarantee for Bayesian methods with heavy-tailed losses. In this paper, we will bridge this gap for the Bayesian framework.

2. Mathematical framework

Let (\mathcal{X}, ζ) be a measurable space and $Z = (X, Y)$ be a random variable taking values in $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ with a probability distribution P where $\mathcal{Y} \subset \mathbb{R}$. We assume that the hypothesis class \mathcal{H} is a bounded subset of the space of square-integrable functions $L_2(\mathcal{X}, \zeta)$ with the convex hull $\overline{\mathcal{H}}$.

For a prior distribution μ on \mathcal{H} and a set D of n independent and identically distributed samples $\{Z_1, Z_2, \dots, Z_n\}$ of Z , the *posterior randomization measure* given the data over the hypothesis space \mathcal{H} has a density

$$p_D(h) \propto \prod_{i=1}^n Q(Z_i | h) = \exp \left\{ - \sum_{i=1}^n \ell(Z_i, h) \right\}$$

with respect to μ , where $\prod_{i=1}^n Q(Z_i | h)$ is called generalized likelihood function and $\ell : \mathcal{Z} \times \overline{\mathcal{H}} \rightarrow \mathbb{R}$ is a function defined by $\ell(Z, h) = -\log Q(Z | h)$, hereafter referred to as the *loss function*.

In the standard Bayesian setting, $Q(Z | h) = P(Z | h)$ where $P(Z | h)$ is the regular density function and the posterior randomization measure is just the posterior distribution. When $Q(Z | h) \neq P(Z | h)$, this setting becomes the quasi-Bayesian approach. For various problems of Bayesian learning using *mean-field variational inference*, the generalized likelihood function belongs to a family of functions that can reasonably approximate the likelihood function. In the “safe Bayesian” framework [4], $Q(Z | h) = [P(Z | h)]^\eta$ where η is a tuned parameter obtained by minimizing a cumulative log-loss. It is worth noticing

that there is a connection between Bayesian setting and PAC-Bayes which has been discussed elsewhere [see e.g. 31, and the references therein].

For a given set of samples D , the *generalized Bayes estimator* is defined as

$$\hat{h} = \int_{\mathcal{H}} p_D(h) h d\mu.$$

Predictions with generalized Bayes estimator are obtained by taking the average of the prediction of the hypotheses in h . The estimator, thus, does not necessarily belong to \mathcal{H} and is an improper estimator. This type of estimator has appeared in various contexts in machine learning. For example, as noted in [32], the safe Bayesian algorithm can be regarded as just running the standard Hedge-algorithm [33] and then making a Cesaro-averaged prediction of the previous Hedge predictions. Similarly, the Weighted Average algorithm [34, 35] makes prediction based on the weighted average predictions of all the hypotheses in the hypothesis space with the weight function

$$w(h) = \exp\left(-c_1 \sum_{i=1}^n |h(X_i) - Y_i|^{c_2}\right)$$

and thus fits into this framework.

We define the *risk function* as $R(h) = \mathbb{E}_{Z \sim P}[\ell(Z, h)]$ and the set of hypotheses whose risks are less than or equal to a threshold value γ as $\mathcal{H}_\gamma = \{h \in \mathcal{H} : R(h) \leq \gamma\}$. For convenience, we assume that

$$\inf_{h \in \mathcal{H}} R(h) = \inf\{\gamma : \mu(\mathcal{H}_\gamma) > 0\} := \gamma^*. \quad (1)$$

Here, γ^* can be considered as “optimal risk”.

We note that this assumption can be relaxed because the set $\{h \in \mathcal{H} : R(h) < \gamma^*\}$ has measure 0. The rationale is that a single best hypothesis is meaningless in the Bayesian setting when \mathcal{H} is uncountable. Hence, we should compare the generalized Bayes estimator to a set of good hypotheses that has a positive measure, as suggested in [36, 35]. The measure of such a set of “good hypotheses” plays a central role in our analyses and directly influences the concentration rates.

In this paper, we are interested in the concentration of the posterior around the set of optimal hypotheses \mathcal{H}_{γ^*} and the convergence properties of the generalized Bayes estimator \hat{h} . Our mathematical framework is designed to analyze the problem of Bayesian learning for unbounded and/or heavy tail losses. We recall that a random variable S is said to have a heavy right tail distribution if

$$\lim_{s \rightarrow \infty} e^{\lambda s} \mathbb{P}[S > s] = \infty$$

for all $\lambda > 0$ and the definition is similar for a heavy left tail distribution. Learning with a heavy-tailed loss means that $\ell(Z, h)$ has a heavy tail distribution from some or all hypotheses $h \in \mathcal{H}$. To enable analyses of fast concentration rates, we impose the following regularity conditions:

Assumption 1 (Regularity condition for risk function). *The risk function R is convex and Lipschitz on $\overline{\mathcal{H}}$.*

We observe that although the risk function R is convex on the convex hull $\overline{\mathcal{H}}$, it may still have multiple global minimizers on \mathcal{H} because we do not put any additional assumption on the geometry of \mathcal{H} . Figure 1 gives an example where this scenario happens. In this example, the convex function $f(x, y) = x^2 + y^2$ achieves the global minimum at two different points $(-1, 0)$ and $(1, 0)$.

Assumption 2 (Multi-scale Bernstein's condition). *There exist a finite partition of $\mathcal{H} = \cup_{i \in I} \mathcal{H}_i$, positive constants $B = \{B_i\}_{i \in I}$, constants $\alpha = \{\alpha_i\}_{i \in I}$ in $(0, 1]$, and a finite set $\mathcal{H}^* = \{h_i^*\}_{i \in I} \subset \mathcal{H}_{\gamma^*}$ such that*

$$\mathbb{E}[\ell(Z, h) - \ell(Z, h_i^*)]^2 \leq B_i [R(h) - \gamma^*]^{\alpha_i}$$

for all $i \in I$ and $h \in \mathcal{H}_i$.

The multi-scale Bernstein's condition is a generalization of the classical Bernstein's condition introduced in [10] to analyze fast convergence rates of the empirical risk minimizer estimator in unbounded losses settings. If a loss function satisfies the Bernstein's condition, then it also satisfies the multi-scale Bernstein's condition. However, while the Bernstein's condition forces the risk

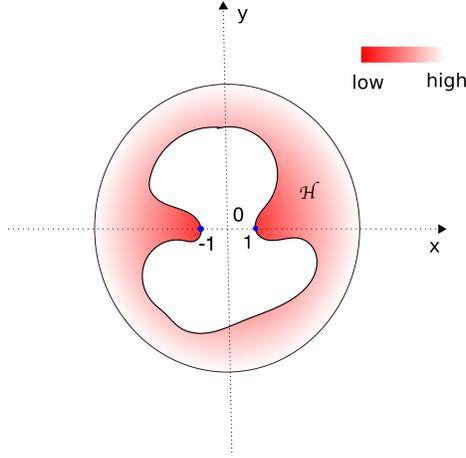


Figure 1: An example where a convex function $f(x, y) = x^2 + y^2$ has two global minimizers (blue dots) on a non-convex set \mathcal{H} . The heat map represents the value of $f(x, y)$ on \mathcal{H} .

function to have a unique minimizer [9], the multi-scale Bernstein's condition does not have this restriction. Note that this condition implies that \mathcal{H}_{γ^*} is not empty.

Assumption 3 (Regularity condition for prior). *There exist $C_1(\mu) > 0$, $C_2(\mu) > 0$, and $\kappa > 0$ such that*

$$\mu(B(h^*, \epsilon)) \geq C_1 \exp(-C_2 \epsilon^{-\kappa})$$

for all $h^* \in \mathcal{H}_{\gamma^*}$. Here, $B(h^*, \epsilon)$ is the ball in $L_2(\mathcal{X}, \zeta)$ with the center h^* and the radius ϵ .

Assumption 3 belongs to a class of regularity assumption called *prior mass assumption* and requires that the prior measures put a sufficient amount of mass near \mathcal{H}_{γ^*} . Such an assumption is standard in the analysis of the convergence rate of posterior measures and can be verified for a broad class of probability distributions [37]. For example, this condition holds for the uniform distribution and the truncated normal distribution on any compact finite-dimensional manifold.

We also need to impose some conditions on the complexity of the hypothesis space. For convenience, let \mathcal{G} denote the set of all functions $g : \mathcal{Z} \rightarrow \mathbb{R}$ such that $g(Z) = \ell(Z, h)$ for some $h \in \mathcal{H}$. For $\epsilon > 0$, let $\mathcal{N}(\epsilon, \mathcal{G}, L_2(P))$ be the *covering number* of $(\mathcal{G}, L_2(P))$; that is, $\mathcal{N}(\epsilon, \mathcal{G}, L_2(P))$ is the minimal number of balls of radius ϵ needed to cover \mathcal{G} . We define the *universal metric entropy* of \mathcal{G} by

$$H(\epsilon, \mathcal{G}) = \sup_Q \log \mathcal{N}(\epsilon, \mathcal{G}, L_2(Q)),$$

where the supremum is taken over the set of all probability measures Q concentrated on some finite subset of \mathcal{Z} . We make the following two assumptions regarding the complexity of \mathcal{G} .

Assumption 4 (Finite covering number). *There exist $\mathcal{C}_1 \geq 1$ and $K_1 \geq 1$ such that*

$$\log \mathcal{N}(\epsilon, \mathcal{G}, L_2(P)) \leq \mathcal{C}_1 \log(K_1/\epsilon) \quad \forall \epsilon \in (0, K_1].$$

Assumption 5 (Universal entropy bounds). *There exist $\mathcal{C}_2 \geq 1$ and $K_2 \geq 1$ such that*

$$H(\epsilon, \mathcal{G}) \leq \mathcal{C}_2 \log(K_2/\epsilon) \quad \forall \epsilon \in (0, K_2].$$

Denote $\mathcal{C} = \max\{\mathcal{C}_1, \mathcal{C}_2\}$. From now on, we will use \mathcal{C} as the common constant for both Assumptions 4 and 5.

Finally, we need a way to control the heavy tails of the loss functions. We employ the integrability condition of the envelope function, which has been studied previously in [10, 13].

Assumption 6 (Integrability of the envelope function). *There exist $W > 0$ and $r \geq 4\mathcal{C}$ such that*

$$\left(\mathbb{E} \sup_{g \in \mathcal{G}} |g|^r \right)^{1/r} \leq W.$$

For convenience, we denote the losses $\ell(Z, h)$ by $\ell(h)$, and define the empirical loss:

$$\ell_D(h) = \frac{1}{n} \sum_{i=1}^n \ell(Z_i, h).$$

For each hypothesis $h_0 \in \mathcal{H}$, we define a ball $\mathcal{B}(h_0, \epsilon)$ of radius $\epsilon > 0$ such that

$$\mathcal{B}(h_0, \epsilon) = \left\{ h \in \mathcal{H} : \left\{ \mathbb{E}[\ell(h) - \ell(h_0)]^2 \right\}^{1/2} \leq \epsilon \right\}.$$

It is worth noticing that if a hypothesis belongs to a ball $\mathcal{B}(h^*, \epsilon)$ for some $h^* \in \mathcal{H}_{\gamma^*}$, then its risk is bounded by $\gamma^* + \epsilon$. To be specific,

$$\bigcup_{h^* \in \mathcal{H}_{\gamma^*}} \mathcal{B}(h^*, \epsilon) \subset \mathcal{H}_{\gamma^* + \epsilon}.$$

3. Fast concentration rates

In this section, we prove that generalized Bayes estimators achieve fast rate learning under the Assumptions introduced in the previous section. Our proof contains two main steps:

Step 1: We prove that the posterior distribution concentrates around the set of optimal hypotheses \mathcal{H}_{γ^*} exponentially fast as n goes to infinity. In other words, we prove that the posterior distribution of the hypotheses which are far away from \mathcal{H}_{γ^*} converges to 0 exponentially fast. The main technique of this step is proving that the difference in the empirical loss between a hypothesis which is close to \mathcal{H}_{γ^*} and a hypothesis which is far away from \mathcal{H}_{γ^*} is sufficiently large.

Step 2: We bound the convergence rate of the generalized Bayes estimator to the optimal risk γ^* . To do so, we show that the generalized Bayes estimator is very close to the average of all hypotheses near \mathcal{H}_{γ^*} . This is due to the fact that the posterior distribution of the hypotheses far away from \mathcal{H}_{γ^*} is small, which was proved in Step 1. Therefore, the risk of the generalized Bayes estimator is close to the average risk of all hypotheses near \mathcal{H}_{γ^*} , which is also close to γ^* .

In the rest of the paper, for some $\beta > 0$, let $\epsilon = n^{-\beta}$ and \mathcal{H}^ϵ denote the finite set containing \mathcal{H}^* such that

$$\bigcup_{h \in \mathcal{H}^\epsilon} \mathcal{B}(h, \epsilon) = \mathcal{H}$$

and that $|\mathcal{H}^\epsilon| \leq (K/\epsilon)^C + |\mathcal{H}^*|$, where \mathcal{H}^* is defined in Assumption 2. Note that Assumption 4 guarantees the existence of \mathcal{H}^ϵ . We now provide the details for the proof of these two steps.

3.1. Posterior concentration

Theorem 1 (Posterior concentration). *Assume that Assumptions 1 – 6 hold.*

Let β be a positive number such that

$$\beta < \max \left\{ \frac{1 - 2\sqrt{\mathcal{C}/r}}{2 - \min_{i \in I} \alpha_i}, \frac{1}{1 + \kappa} \right\}.$$

Then, for any $\delta \in (0, 1)$, there exist $C_{r,\beta}, C'_{r,\beta} > 0$ and $N_{\delta,r,B,\alpha,\kappa} > 0$ such that for $n \geq N_{\delta,r,B,\alpha,\kappa}$ and $\epsilon = n^{-\beta}$, we have:

$$\sup_{h \in \mathcal{H} \setminus \mathcal{H}_{\gamma^* + C_{\delta,r,\beta}\epsilon}} p_D(h) \leq \frac{1}{C_1} \exp \left\{ -\frac{1}{2} \left[C_{r,\beta} + \left(\frac{C'_{r,\beta}}{\delta} \right)^{1/[2\sqrt{\mathcal{C}r}]} \right] n^{1-\beta} \right\}$$

with probability at least $1 - \delta$.

The detailed proof of this theorem is provided in the Appendix. Here, we want to give some insights about the proof's arguments. Let us consider the simplest case when \mathcal{H} is finite and the optimal hypothesis h^* is unique. In this setting, for any other hypothesis $h \in \mathcal{H}$, by the strong law of large numbers, we have

$$\ell_D(h) - \ell_D(h^*) \approx R(h) - R(h^*)$$

as the sample size n goes to infinity. Informally, this implies that $p_D(h)/p_D(h^*) \rightarrow 0$. Since \mathcal{H} is finite and $\sum_{h \in \mathcal{H}} p_D(h) = 1$, we deduce that the distribution concentrates around $\mathcal{H}_{\gamma^*} = \{h^*\}$.

In the case when \mathcal{H} is infinite, comparing between two hypotheses becomes less meaningful. To extend the result, we need to provide a uniform bound on

$\ell_D(h) - \ell_D(h')$ for all $h \in U_1$ and $h' \in U_2$, where U_2 is a neighborhood of h^* and U_1 is a set that covers most of the outside of U_2 . This estimate is obtained by a combination of the following two Lemmas, of which the optimal hypothesis h^* acts as an intermediary for comparisons.

Lemma 1. *Assume that Assumptions 1, 2, 4, 5, and 6 hold. For any $\beta < 1 - 2\sqrt{\mathcal{C}/r}$, there exists $C_{r,\beta}, C'_{r,\beta} > 0$ and such that for all $n \in \mathbb{N}$ and $\delta \in (0, 1)$, we have:*

$$|\ell_D(h) - \ell_D(h_0)| \leq \left[C_{r,\beta} + \left(\frac{C'_{r,\beta}}{\delta} \right)^{1/[2\sqrt{\mathcal{C}r}]} \right] \epsilon, \quad \forall h_0 \in \mathcal{H}^\epsilon, h \in \mathcal{B}(h_0, \epsilon)$$

with probability at least $1 - \delta$.

Lemma 2. *Assume that Assumptions 1, 2, 4, and 6 hold. For any $a > 0$, $\delta \in (0, 1)$, and a positive number β satisfying*

$$\beta < (1 - 2\sqrt{\mathcal{C}/r}) / (2 - \alpha_i) \quad \forall i \in I,$$

there exists $N_{a,\delta,r,B,\alpha} > 0$ such that for $n \geq N_{a,\delta,r,B,\alpha}$, we have

$$\forall h \in \mathcal{H}^\epsilon \setminus \mathcal{H}_{\gamma^*+a\epsilon}, \exists h^* \in \mathcal{H}^* : \ell_D(h) - \ell_D(h^*) > \frac{a\epsilon}{4}$$

with probability at least $1 - \delta$. Here, I and \mathcal{H}^* are defined in Assumption 2.

Lemma 1 is a consequence of Lemma 3.5 in [10] and Lemma 2 is Theorem 3.2 in [10]. Lemma 2 ensures that a hypothesis that has small empirical loss will also have small risk. This result provides an alternative to concentration bound, which may not exist. It is similar to the techniques of using one-sided inequalities for learning without concentration bound, established in [38].

Finally, when the optimal hypothesis is not unique, we need to utilize the multi-scale Bernstein's condition to partition the hypothesis spaces into regions where local behavior of the empirical loss function can be controlled, and combine the estimates in later steps. We note that the feasibility of this approach comes from the fact that the multi-scale Bernstein's condition is a local condition.

From now on, to ease the notation, we denote

$$C_{\delta,r,\beta} = \frac{1}{2} \left[C_{r,\beta} + \left(\frac{C'_{r,\beta}}{\delta} \right)^{1/[2\sqrt{\mathcal{C}r}]} \right].$$

3.2. Learning rates

Theorem 2 (Learning rate). *Assume that Assumptions 1 – 6 hold. Let β be a positive number satisfying*

$$\beta < \max \left\{ \frac{1 - 2\sqrt{\mathcal{C}/r}}{2 - \min_{i \in I} \alpha_i}, \frac{1}{1 + \kappa} \right\}.$$

Then, for any $\delta \in (0, 1)$, there exists $N_{\delta, r, \beta, \mu, \kappa} > 0$ such that

$$\mathbb{P} \left(\hat{h}_n \in \mathcal{H}_{\gamma^* + 2C_{\delta, r, \beta, \mu, \kappa} \epsilon} \right) \geq 1 - \delta,$$

for all $n \geq N_{\delta, r, \beta, \mu, \kappa}$ and $\epsilon = n^{-\beta}$.

Proof. We define

$$M = \sup_{h \in \mathcal{H}} \|h\|_2 < \infty, \quad \text{and} \quad \nu = \int_{\mathcal{H}_{\gamma^* + C_{\delta, r, \beta, \mu, \kappa} \epsilon}} p_D(h) d\mu \leq 1.$$

Note that M is finite because \mathcal{H} is a bounded subset of $L^2(\mathcal{X}, \zeta)$. On the other hand, by Theorem 1, with probability at least $1 - \delta$:

$$1 - \nu = \int_{\mathcal{H} \setminus \mathcal{H}_{\gamma^* + C_{\delta, r, \beta, \mu, \kappa} \epsilon}} p_D(h) d\mu \leq \frac{\exp\{-C_{\delta, r, \beta, \mu, \kappa} n \epsilon\}}{C_1}.$$

Hence, when n is sufficient large, we have $\nu > 0$ with probability at least $1 - \delta$.

By Assumption 1, R is convex and Lipchitz in $\overline{\mathcal{H}}$. Therefore,

$$\int_{\mathcal{H}} h p_D(h) d\mu \quad \text{and} \quad \int_{\mathcal{H}_{\gamma^* + C_{\delta, r, \beta, \mu, \kappa} \epsilon}} h \frac{p_D(h)}{\nu} d\mu$$

belong to $\overline{\mathcal{H}}$ and there exists a Lipchitz constant L such that

$$\begin{aligned} & \left| R \left(\int_{\mathcal{H}} h p_D(h) d\mu \right) - R \left(\int_{\mathcal{H}_{\gamma^* + C_{\delta, r, \beta, \mu, \kappa} \epsilon}} h \frac{p_D(h)}{\nu} d\mu \right) \right| \\ & \leq L \left\| \int_{\mathcal{H}} h p_D(h) d\mu - \int_{\mathcal{H}_{\gamma^* + C_{\delta, r, \beta, \mu, \kappa} \epsilon}} h \frac{p_D(h)}{\nu} d\mu \right\|_2. \end{aligned}$$

We deduce that

$$\begin{aligned} R(\hat{h}_n) &= R \left(\int_{\mathcal{H}} h p_D(h) d\mu \right) \leq R \left(\int_{\mathcal{H}_{\gamma^* + C_{\delta, r, \beta, \mu, \kappa} \epsilon}} h \frac{p_D(h)}{\nu} d\mu \right) \\ &+ L \left\| \int_{\mathcal{H} \setminus \mathcal{H}_{\gamma^* + C_{\delta, r, \beta, \mu, \kappa} \epsilon}} h p_D(h) d\mu \right\|_2 + \frac{1 - \nu}{\nu} L \left\| \int_{\mathcal{H}_{\gamma^* + C_{\delta, r, \beta, \mu, \kappa} \epsilon}} h p_D(h) d\mu \right\|_2. \end{aligned}$$

We have

$$\begin{aligned} R\left(\int_{\mathcal{H}_{\gamma^*+C_{\delta,r,\beta\epsilon}}} h \frac{p_D(h)}{\nu} d\mu\right) &\leq \int_{\mathcal{H}_{\gamma^*+C_{\delta,r,\beta\epsilon}}} R(h) \frac{p_D(h)}{\nu} d\mu \\ &\leq \int_{\mathcal{H}_{\gamma^*+C_{\delta,r,\beta\epsilon}}} (\gamma^* + C_{\delta,r,\beta\epsilon}) \frac{p_D(h)}{\nu} d\mu = \gamma^* + C_{\delta,r,\beta\epsilon}. \end{aligned}$$

Moreover,

$$\left\| \int_{\mathcal{H} \setminus \mathcal{H}_{\gamma^*+C_{\delta,r,\beta\epsilon}}} h p_D(h) d\mu \right\|_2 \leq \int_{\mathcal{H} \setminus \mathcal{H}_{\gamma^*+C_{\delta,r,\beta\epsilon}}} \|h\|_2 p_D(h) d\mu \leq M(1-\nu).$$

and

$$\left\| \int_{\mathcal{H}_{\gamma^*+C_{\delta,r,\beta\epsilon}}} h p_D(h) d\mu \right\|_2 \leq \int_{\mathcal{H}_{\gamma^*+C_{\delta,r,\beta\epsilon}}} \|h\|_2 p_D(h) d\mu \leq M\nu.$$

We conclude that with probability at least $1 - \delta$,

$$R(\hat{h}_n) \leq \gamma^* + C_{\delta,r,\beta\epsilon} + 2LM(1-\nu) \leq \gamma^* + C_{\delta,r,\beta\epsilon} + 2LM \frac{\exp\{-C_{\delta,r,\beta}n\epsilon\}}{C_1}.$$

Hence, when n is sufficiently large, we have

$$R(\hat{h}) \leq \gamma^* + 2C_{\delta,r,\beta\epsilon}$$

with probability at least $1 - \delta$, which completes the proof for the theorem. \square

The result of Theorem 2 implies

Corollary 1. *For all $\delta \in (0, 1)$, $R(\hat{h}_n) = \gamma^* + \mathcal{O}(n^{-\beta})$ with probability at least $1 - \delta$, where*

$$\beta < \max \left\{ \frac{1 - 2\sqrt{\mathcal{C}/r}}{2 - \min_{i \in I} \alpha_i}, \frac{1}{1 + \kappa} \right\}.$$

When r is sufficiently large, $\min_{i \in I} \alpha_i = 1$, and κ is sufficiently small, we achieve convergence rates arbitrarily close to $\mathcal{O}(n^{-1})$.

Hence, fast learning rates for generalized Bayesian estimators are available within our framework. In general, the order of convergence depends on the regularity of the loss function (via the multi-scale Bernstein's order) and the balance between the complexity of the hypothesis class and the thickness of the tail of the loss's distribution.

4. Robustness of Bayesian linear regression

In this section, we will apply our results to show that Bayesian linear regression is robust to heavy-tailed distributions. To be specific, we consider the following standard linear regression setting:

$$Y_i = \mathbf{X}_i u_0 + \epsilon_i, \quad i = 1, 2, \dots, n$$

where $Y_i \in \mathbb{R}$, $\mathbf{X}_i \in \mathcal{X} \in \mathbb{R}^d$, $u_0 \in \mathbb{R}^d$, and ϵ_i are i.i.d random variables which follow a t-distribution with degree of freedom k .

We will prove that even if we do not know that ϵ_i is heavy-tailed and just assume that ϵ_i follows a standard normal distribution, the Bayesian linear regression still achieves fast rate learning. Given a proper prior π_u for u , the posterior distribution of u has the following form:

$$p_D(u \mid \{Y_i\}_{i=1}^n) \propto \pi_u \exp \left\{ -\sum_{i=1}^n \frac{(Y_i - \mathbf{X}_i u)^2}{2} \right\}, \quad (2)$$

which corresponds to our setting with the loss function $\ell(Y, \mathbf{X}, u) = (Y - \mathbf{X}u)^2$.

Theorem 3. *Assume that $\|u_0\|_2 \leq M_u$, \mathcal{X} is bounded in $\|\cdot\|_2$ by M_X , $k > 4d$, and π_u is regular (Assumption 3). Let β be a positive number satisfying*

$$\beta < \min\{1 - 2\sqrt{d/k}, 1/(1 + \kappa)\}.$$

The Bayesian linear regression estimator

$$\hat{u} = \int u \cdot p_D(u \mid \{Y_i\}_{i=1}^n) du$$

achieves learning rate $n^{-\beta}$.

The proof of this theorem is in the Appendix. We observe that when $k > 16d$, Theorem 3 implies the Bayesian linear regression estimator achieves fast learning rate. It is worth noticing that most of the common priors on a bounded set of \mathbb{R}^d (for example, uniform distribution) satisfy the regularity condition for prior (Assumption 3) with any $\kappa > 0$.

Simulations. To illustrate the result, we use the R-platform to simulate data from the following model:

$$Y_i = 1 + X_i^{(1)} + X_i^{(2)} + \epsilon_i, \quad i = 1, 2, \dots, n$$

where $\{X_i^{(1)}\}, \{X_i^{(2)}\}$ are i.i.d. random variables that follow a truncated standard normal distribution (the truncation value is 1), and ϵ_i are i.i.d random variables which follow a t-distribution with degree of freedom $k = 5, 10, 20$. For each degree of freedom, we vary the sample size from 10 to 10240 ($n = 10, 20, 40, 80, 160, 320, 640, 1280, 2560, 5120, 10240$), and for each sample size, we simulate 100 data sets. We analyze each data set using the standard linear regression (ERM) and the Bayesian linear regression (2) with a uniform prior on the ball which is centered at 0 and has radius 10. We explore the generalized posterior distribution of the coefficients using the Metropolis algorithm implemented in the R function `MCMCmetrop1R` from the package `MCMCpack` [39]. We discard the first 20000 iterations of the Markov chain Monte Carlo and use the next 100000 iterations to approximate the Bayesian linear regression estimator. Then, we apply Monte Carlo method to approximate the risk of the estimators and fit a linear regression between the risk (in log-scale) and the sample size (in log-scale) to approximate the rate of convergence of the ERM and the Bayesian linear regression (e.g. Figure 2). We summarize the result of our simulations in Table 1. The result confirms that Bayesian linear regression is robust to heavy-tailed noises. We note that the empirical convergence rate of the Bayesian linear regression (as well as the ERM, which has been investigated in [10]) is faster than its theoretical bound in Theorem 3.

Degree of freedom	ERM	Bayesian
5	-0.984	-0.951
10	-1.001	-0.966
20	-1.046	-0.996

Table 1: The approximated rate of convergence of the two estimators with $k = 5, 10, 20$.

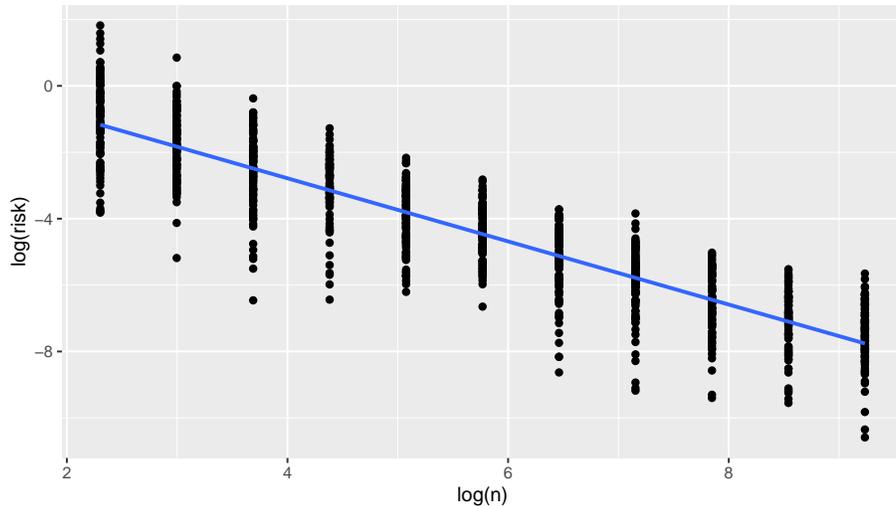


Figure 2: A linear regression between the risk (in log-scale) of the Bayesian linear regression and the sample size (in log-scale) when the errors follow the t -distribution with 5 degrees of freedom. The slope of the fitted line approximates the rate of convergence of the estimator.

5. Discussions and Conclusions

The result of this paper indicates that learning with Bayesian estimators and heavy-tailed losses can obtain convergence rates up to an essential order

$$\mathcal{O}\left(n^{-(1-2\sqrt{C/r})/(2-\min\{\alpha\})}\right)$$

where α is the multi-scale Bernstein's order and r is the degree of integrability of the loss. This result is consistent with previous works using a frequentist approach [10]. We note that for bounded and strongly convex losses, our assumptions can be validated with $\alpha = 1$, $I = 1$, and $r = \infty$ and this reduces to the convergence rate $\mathcal{O}(1/n)$.

There are several avenues for improvement. Firstly, in this work, we consider a setting where the generalized likelihood function has the form $\prod_{i=1}^n Q(Z_i | h)$. In some scenarios, for example, when data are dependent, this setting may not hold. It would be interesting to see if concentration and learning rates retain in those cases. Secondly, although our framework (which relies on the multi-scale

Bernstein’s condition) allows us to analyze the convergence of generalized Bayes estimators in more general settings than previous approaches, our result requires high-order moments of the loss to guarantee convergence. Recently, there has been a growing interest in fast learning rate for convex losses using the small-ball condition [15], which requires only low-order moments. We would like to extend the result in this paper to study and adapt this condition to the case when the optimal hypothesis is non-unique.

Finally, the simulations confirm the robustness of Bayesian linear regression to heavy-tailed noises. This is an assurance for end-users that Bayesian linear regression is not vulnerable to the violation of the assumption of normal errors. In particular, no special treatment is needed when the errors follow a t-distribution and the Bayesian estimates converge to the true values at the same rate as their frequentist counterparts. It is worth noticing that the simulations indicate that the convergence rate is $\mathcal{O}(1/n)$, which means that our theoretical upper bounds may not be optimal. An interesting direction for future research is to derive sharper upper bounds and/or lower bounds.

Acknowledgments

LSTH was supported by startup funds from Dalhousie University, the Canada Research Chairs program, the NSERC Discovery Grant RGPIN-2018-05447, and the NSERC Discovery Launch Supplement DGEER-2018-00181.

References

- [1] J. Rousseau, K. Mengersen, Asymptotic behaviour of the posterior distribution in overfitted mixture models, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 73 (5) (2011) 689–710.
- [2] J. Arbel, G. Gayraud, J. Rousseau, Bayesian optimal adaptive estimation using a sieve prior, *Scandinavian journal of statistics* 40 (3) (2013) 549–570.
- [3] J. Rousseau, On the frequentist properties of Bayesian nonparametric methods, *Annual Review of Statistics and Its Application* 3 (2016) 211–231.
- [4] P. Grünwald, J. S. Jones, J. de Winter, É. Smith, Safe Learning: bridging the gap between Bayes, MDL and statistical learning theory via empirical convexity., in: *COLT*, 2011, pp. 397–420.
- [5] J.-Y. Audibert, A. B. Tsybakov, Fast learning rates for plug-in classifiers, *The Annals of statistics* 35 (2) (2007) 608–633.
- [6] V. Dinh, L. S. T. Ho, N. V. Cuong, D. Nguyen, B. T. Nguyen, Learning from non-iid data: Fast rates for the one-vs-all multiclass plug-in classifiers, in: *Theory and Applications of Models of Computation*, Springer, 2015, pp. 375–387.
- [7] N. A. Mehta, R. C. Williamson, From stochastic mixability to fast rates, in: *Advances in Neural Information Processing Systems*, 2014, pp. 1197–1205.
- [8] P. L. Bartlett, S. Mendelson, Empirical minimization, *Probability Theory and Related Fields* 135 (3) (2006) 311–334.
- [9] T. van Erven, P. D. Grünwald, N. A. Mehta, M. D. Reid, R. C. Williamson, Fast rates in statistical and online learning, *Journal of Machine Learning Research* 16 (2015) 1793–1861.
- [10] V. C. Dinh, L. S. Ho, B. Nguyen, D. Nguyen, Fast learning rates with heavy-tailed losses, in: *Advances in Neural Information Processing Systems*, 2016, pp. 505–513.

- [11] T. Zhang, From ϵ -entropy to KL-entropy: Analysis of minimum information complexity density estimation, *The Annals of Statistics* 34 (5) (2006) 2180–2210.
- [12] T. Zhang, Information-theoretic upper and lower bounds for statistical estimation, *IEEE Transactions on Information Theory* 52 (4) (2006) 1307–1321.
- [13] G. Lecué, S. Mendelson, General nonexact oracle inequalities for classes with a subexponential envelope, *The Annals of Statistics* 40 (2) (2012) 832–860.
- [14] D. Hsu, S. Sabato, Loss minimization and parameter estimation with heavy tails, *Journal of Machine Learning Research* 17 (18) (2016) 1–40.
- [15] S. Mendelson, On aggregation for heavy-tailed classes, *Probability Theory and Related Fields* (2017) 1–34.
- [16] P. D. Grünwald, N. A. Mehta, Fast rates with unbounded losses, arXiv preprint arXiv:1605.00252.
- [17] F. E. Bachl, F. Lindgren, D. L. Borchers, J. B. Illian, inlabru: an R package for Bayesian spatial modelling from ecological survey data, *Methods in Ecology and Evolution* 10 (6) (2019) 760–766.
- [18] M. S. Gill, L. S. T. Ho, G. Baele, P. Lemey, M. A. Suchard, A relaxed directional random walk model for phylogenetic trait evolution, *Systematic biology* 66 (3) (2017) 299–319.
- [19] L. S. T. Ho, J. Xu, F. W. Crawford, V. N. Minin, M. A. Suchard, Birth/birth-death processes and their computable transition probabilities with biological applications, *Journal of mathematical biology* 76 (4) (2018) 911–944.
- [20] L. S. T. Ho, F. W. Crawford, M. A. Suchard, et al., Direct likelihood-based inference for discretely observed stochastic compartmental models of infectious disease, *The Annals of Applied Statistics* 12 (3) (2018) 1993–2021.

- [21] J. Geweke, G. Gowrisankaran, R. J. Town, Bayesian inference for hospital quality in a selection model, *Econometrica* 71 (4) (2003) 1215–1238.
- [22] C. Brownlees, E. Joly, G. Lugosi, Empirical risk minimization for heavy-tailed losses, *The Annals of Statistics* 43 (6) (2015) 2507–2536.
- [23] G. Lugosi, S. Mendelson, Mean estimation and regression under heavy-tailed distributions: A survey, *Foundations of Computational Mathematics* 19 (5) (2019) 1145–1190.
- [24] O. Bachem, M. Lucic, S. H. Hassani, A. Krause, Uniform deviation bounds for k-means clustering, in: *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, JMLR. org, 2017, pp. 283–291.
- [25] A. Christmann, I. Steinwart, A. van Messem, On consistency and robustness properties of support vector machines for heavy-tailed distributions, *Statistics and Its Interface* 2 (3) (2009) 311–327.
- [26] Q. Han, J. A. Wellner, et al., Convergence rates of least squares regression estimators with heavy-tailed errors, *The Annals of Statistics* 47 (4) (2019) 2286–2319.
- [27] X. Nguyen, et al., Borrowing strength in hierarchical bayes: Posterior concentration of the Dirichlet base measure, *Bernoulli* 22 (3) (2016) 1535–1571.
- [28] V. Dinh, A. E. Rundell, G. T. Buzzard, Convergence of gridy gibbs sampling and other perturbed markov chains, *Journal of Statistical Computation and Simulation* 87 (7) (2017) 1379–1400.
- [29] P. Grünwald, T. Van Ommen, et al., Inconsistency of Bayesian inference for misspecified linear models, and a proposal for repairing it, *Bayesian Analysis* 12 (4) (2017) 1069–1103.
- [30] R. de Heide, A. Kirichenko, N. Mehta, P. Grünwald, Safe-Bayesian Generalized Linear Regression, *arXiv preprint arXiv:1910.09227*.

- [31] P. Germain, F. Bach, A. Lacoste, S. Lacoste-Julien, PAC-Bayesian theory meets Bayesian inference, in: *Advances in Neural Information Processing Systems*, 2016, pp. 1884–1892.
- [32] P. Grünwald, The safe Bayesian, in: *International Conference on Algorithmic Learning Theory*, Springer, 2012, pp. 169–183.
- [33] Y. Freund, R. E. Schapire, A decision-theoretic generalization of on-line learning and an application to boosting, in: *European conference on computational learning theory*, Springer, 1995, pp. 23–37.
- [34] J. Kivinen, M. K. Warmuth, Averaging expert predictions, in: *European Conference on Computational Learning Theory*, Springer, 1999, pp. 153–167.
- [35] N. V. Cuong, L. S. T. Ho, V. Dinh, Generalization and robustness of batched weighted average algorithm with V-geometrically ergodic Markov data, in: *International Conference on Algorithmic Learning Theory*, Springer, 2013, pp. 264–278.
- [36] Y. Freund, Y. Mansour, R. E. Schapire, Generalization bounds for averaged classifiers, *Annals of Statistics* (2004) 1698–1722.
- [37] S. Ghosal, J. K. Ghosh, A. W. Van Der Vaart, Convergence rates of posterior distributions, *Annals of Statistics* (2000) 500–531.
- [38] S. Mendelson, Learning without concentration, *Journal of the ACM (JACM)* 62 (3) (2015) 21.
- [39] A. D. Martin, K. M. Quinn, J. H. Park, MCMCpack: Markov chain Monte Carlo in R, *Journal of Statistical Software* 42 (9) (2011) 22.
URL <http://www.jstatsoft.org/v42/i09/>

Appendix A. Detailed proofs

Proof of Theorem 1. We denote

$$r_n(h) = \exp\{-\ell_D(h)\}.$$

Then, the posterior can be calculated by the following formula:

$$p_D(h) = \left[\frac{r_n(h)}{\|r_n(h)\|_n} \right]^n,$$

where

$$\|r_n(h)\|_n = \left(\int_{\mathcal{H}} |r_n(h)|^n d\mu \right)^{1/n}.$$

For any $i \in I$, we apply Lemma 2 with $a = 12C_{\delta,r,\beta}$ to obtain

$$\ell_D(h) - \ell_D(h_i^*) > 3C_{\delta,r,\beta}\epsilon, \quad \forall h \in (\mathcal{H}_i \setminus \mathcal{H}_{\gamma^*+12C_{\delta,r,\beta}\epsilon}) \cap \mathcal{H}^\epsilon$$

with probability $1 - \delta$.

By Lemma 1, we derive that

$$\ell_D(h) - \ell_D(h') > C_{\delta,r,\beta}\epsilon, \quad \forall h \in \mathcal{H}_i \setminus \mathcal{H}_{\gamma^*+12C_{\delta,r,\beta}\epsilon}, \quad h' \in \mathcal{B}(h_i^*, \epsilon)$$

with probability $1 - 3\delta$.

Hence,

$$r_n(h) \leq e^{-C_{\delta,r,\beta}\epsilon} r_n(h'), \quad \forall h \in \mathcal{H}_i \setminus \mathcal{H}_{\gamma^*+12C_{\delta,r,\beta}\epsilon}, \quad h' \in \mathcal{B}(h_i^*, \epsilon)$$

with probability at least $1 - 3\delta$.

Therefore,

$$\sup_{h \in \mathcal{H}_i \setminus \mathcal{H}_{\gamma^*+12C_{\delta,r,\beta}\epsilon}} r_n(h) \leq e^{-C_{\delta,r,\beta}\epsilon} \inf_{h' \in \mathcal{B}(h_i^*, \epsilon)} r_n(h').$$

with probability at least $1 - 3\delta$.

We have

$$\|r_n\|_n = \left(\int_{\mathcal{H}} |r_n(h)|^n d\mu \right)^{1/n} \geq \left(\int_{\mathcal{B}(h_i^*, \epsilon)} |r_n(h)|^n d\mu \right)^{1/n} = \inf_{h' \in \mathcal{B}(h_i^*, \epsilon)} r_n(h') \mu(\mathcal{B}(h_i^*, \epsilon))^{1/n},$$

with probability at least $1 - 3\delta$.

Consequently, when n is sufficient large,

$$\begin{aligned} \sup_{h \in \mathcal{H}_i \setminus \mathcal{H}_{\gamma^* + 12C_{\delta, r, \beta} \epsilon}} p_D(h) &= \sup_{h \in \mathcal{H}_i \setminus \mathcal{H}_{\gamma^* + 12C_{\delta, r, \beta} \epsilon}} \left(\frac{r_n(h)}{\|r_n\|_n} \right)^n \leq \frac{e^{-2C_{\delta, r, \beta} n \epsilon}}{\mu(\mathcal{B}(h_i^*, \epsilon))} \\ &\leq \frac{1}{C_1} \exp(-n2C_{\delta, r, \beta} \epsilon + C_2 \epsilon^{-\kappa}) \leq \frac{1}{C_1} \exp(-nC_{\delta, r, \beta} \epsilon), \end{aligned}$$

with probability at least $1 - 3\delta$.

Under Assumption 2, I is finite and $\mathcal{H} = \bigcup_{i \in I} \mathcal{H}_i$. Therefore, the proof is completed by taking a union bound over I . \square

Proof of Theorem 3. Let u_0 be the true value of u . We will verify Assumptions 1, 2, 4, and 6. Instead of checking Assumption 5, we will prove Lemma 1 directly.

Assumption 1: The risk function $R(u) = \mathbb{E}[(Y - \mathbf{X}u)^2]$ is convex and Lipschitz in u . Indeed,

$$\begin{aligned} R\left(\frac{u_1 + u_2}{2}\right) &= \mathbb{E}\left[\left(Y - \mathbf{X}\frac{u_1 + u_2}{2}\right)^2\right] \\ &\leq \frac{1}{2}(\mathbb{E}[(Y - \mathbf{X}u_1)^2] + \mathbb{E}[(Y - \mathbf{X}u_2)^2]) \\ &= \frac{1}{2}(R(u_1) + R(u_2)), \end{aligned}$$

and

$$\begin{aligned} |R(u_1) - R(u_2)| &= |\mathbb{E}[(Y - \mathbf{X}u_1)^2 - (Y - \mathbf{X}u_2)^2]| \\ &\leq M_X \|u_2 - u_1\|_2 [2\mathbb{E}|Y - \mathbf{X}u_0| + 4M_X M_u] \\ &\leq M_X \|u_2 - u_1\|_2 \{2[\mathbb{E}(Y - \mathbf{X}u_0)^2]^{1/2} + 4M_X M_u\} \\ &= M_X \left(\frac{2k^{1/2}}{(k-2)^{1/2}} + 4M_X M_u \right) \|u_2 - u_1\|_2. \end{aligned}$$

Assumption 2: We first note that u_0 is the only optimal hypothesis. Indeed,

$$\begin{aligned} R(u) - R(u_0) &= \mathbb{E}[\mathbf{X}(u - u_0)(2Y - \mathbf{X}(u + u_0))] \\ &= \mathbb{E}[\mathbb{E}[\mathbf{X}(u - u_0)(2Y - \mathbf{X}(u + u_0)) \mid \mathbf{X}]] \\ &= \mathbb{E}[\mathbf{X}(u - u_0)(2\mathbb{E}[Y \mid \mathbf{X}] - \mathbf{X}(u + u_0))] \\ &= \mathbb{E}[(\mathbf{X}(u - u_0))^2]. \end{aligned}$$

$$\inf_{u \neq u_0} \frac{R(u) - R(u_0)}{\|u - u_0\|_2^2} = \inf_{\|z\|_2=1} \mathbb{E}([\mathbf{X}z]^2).$$

Since $\mathbb{E}([\mathbf{X}z]^2) > 0$ for all $\|z\|_2 = 1$ then $\inf_{\|z\|_2=1} \mathbb{E}([\mathbf{X}z]^2) \geq D > 0$. Therefore, u_0 is the only optimal hypothesis and $R(u) - R(u_0) \geq D\|u - u_0\|^2$. Note that we have proved $\mathbb{E}\{[(Y - \mathbf{X}u_1)^2 - (Y - \mathbf{X}u_2)^2]^2\} \leq C_0\|u_1 - u_2\|^2$. We conclude that the multi-scale Bernstein condition is satisfied with $\alpha = 1$.

Assumption 4:

$$\begin{aligned} [d_P(u_1, u_2)]^2 &= \mathbb{E}\{[(Y - \mathbf{X}u_1)^2 - (Y - \mathbf{X}u_2)^2]^2\} \\ &= \mathbb{E}\{[\mathbf{X}(u_1 - u_2)]^2 [2Y - \mathbf{X}(u_1 + u_2)]^2\} \\ &\leq M_X \|u_1 - u_2\|_2^2 \mathbb{E}\{[2Y - \mathbf{X}(u_1 + u_2)]^2\} \\ &\leq M_X \|u_1 - u_2\|_2^2 \{8\mathbb{E}[(Y - \mathbf{X}u_0)^2] + 32M_X^2 M_u^2\} \\ &= M_X \left[8\frac{k}{k-2} + 32M_X^2 M_u^2 \right] \|u_1 - u_2\|^2 = C_0 \|u_1 - u_2\|^2. \end{aligned}$$

Therefore, Assumption 4 holds with $C_1 = d$ and $K_1 = M_u$.

Assumption 6:

$$\begin{aligned} \mathbb{E}[\sup_u (Y - \mathbf{X}u)^r] &\leq \mathbb{E}[\sup_u (|Y - \mathbf{X}u_0| + 2M_u M_X)^r] \\ &= \mathbb{E}[(|Y - \mathbf{X}u_0| + 2M_u M_X)^r] \leq W \end{aligned}$$

when $r < k$. Then Assumption 6 is satisfied with any $r \in [4d, k)$.

Lemma 1: Note that

$$\frac{[d_P(u_1, u_2)]^2}{\|u_1 - u_2\|_2^2} = \mathbb{E} \left\{ \left[\mathbf{X} \frac{u_1 - u_2}{\|u_1 - u_2\|_2} \right]^2 [2Y - \mathbf{X}(u_1 + u_2)]^2 \right\} > 0$$

for all u_1, u_2 . Since u_1, u_2 are bounded, we have

$$\frac{[d_P(u_1, u_2)]^2}{\|u_1 - u_2\|_2^2} \geq D > 0.$$

Therefore,

$$\begin{aligned}
& |\ell_D(u_1) - \ell_D(u_2)| \\
&= \left| \frac{1}{n} \sum_{i=1}^n (Y_i - \mathbf{X}_i u_1)^2 - (Y_i - \mathbf{X}_i u_2)^2 \right| \\
&\leq \frac{1}{n} \sum_{i=1}^n 2(|Y_i - \mathbf{X}_i u_0| + M_X M_u) M_X \|u_1 - u_2\|_2 \\
&\leq \frac{1}{n} \sum_{i=1}^n 2(|Y_i - \mathbf{X}_i u_0| + M_X M_u) M_X \frac{d_P(u_1, u_2)}{D}.
\end{aligned}$$

So,

$$\sup_{u_1 \in \mathcal{H}, u_2 \in \mathcal{B}(u_1, \epsilon)} |\ell_D(u_1) - \ell_D(u_2)| \leq \frac{1}{n} \sum_{i=1}^n 2(|Y_i - \mathbf{X}_i u_0| + M_X M_u) M_X \frac{\epsilon}{D}.$$

Note that, for all $M > 0$,

$$\Pr \left(\frac{1}{n} \sum_{i=1}^n |Y_i - \mathbf{X}_i u_0| \geq M \right) \leq \frac{\mathbb{E}|Y_1 - \mathbf{X}_1 u_0|}{M} \leq \frac{[\mathbb{E}(Y_1 - \mathbf{X}_1 u_0)^2]^{1/2}}{M} \leq \frac{k^{1/2}}{M(k-2)^{1/2}}.$$

Hence, we can choose M_δ such that

$$\sup_{u_1 \in \mathcal{H}, u_2 \in \mathcal{B}(u_1, \epsilon)} |\ell_D(u_1) - \ell_D(u_2)| \leq \frac{2(M_\delta + M_X M_u) M_X}{D} \epsilon$$

with probability at least $1 - \delta$.

□