
A Statistical Perspective on Retrieval-Based Models

Soumya Basu^{*1} Ankit Singh Rawat^{*2} Manzil Zaheer^{*3}

Abstract

Many modern high-performing machine learning models increasingly rely on scaling up models, e.g., transformer networks. Simultaneously, a parallel line of work aims to improve the model performance by augmenting an input instance with other (labeled) instances during inference. Examples of such augmentations include task-specific prompts and similar examples retrieved from the training data by a nonparametric component. Despite a growing literature showcasing the promise of these *retrieval-based models*, their theoretical underpinnings remain under-explored. In this paper, we present a formal treatment of retrieval-based models to characterize their performance via a novel statistical perspective. In particular, we study two broad classes of retrieval-based classification approaches: First, we analyze a *local learning* framework that employs an *explicit* local empirical risk minimization based on retrieved examples for each input instance. Interestingly, we show that breaking down the underlying learning task into local sub-tasks enables the model to employ a *low complexity* parametric component to ensure good overall performance. The second class of retrieval-based approaches we explore learns a global model using kernel methods to directly map an input instance and retrieved examples to a prediction, without explicitly solving a local learning task.

1. Introduction

As our world is complex, we need expressive machine learning (ML) models to make high-accuracy predictions on real-world problems. There are multiple ways to increase the expressiveness of an ML model. A popular way is to

^{*}Equal contribution; in alphabetical order ¹Google, Mountain View, USA ²Google Research, New York, USA ³Google DeepMind, New York, USA. Correspondence to: Soumya Basu <basu-soumya@google.com>.

homogeneously scale the size of a parametric model, such as neural networks, which has been behind many recent high-performance models such as GPT-3 (Brown et al., 2020) and ViT (Dosovitskiy et al., 2021). Their performance (accuracy) exhibits a monotonic behavior with increasing model size, as demonstrated by “scaling laws” (Kaplan et al., 2020; Hoffmann et al., 2022). Such large models, however, have their own limitations, including high computation cost, catastrophic forgetting (hard to adapt to changing data), lack of provenance, and poor explainability. Classical instance-based models (Fix & Hodges, 1989), on the other hand, offer many desirable properties by design — efficient data structures, incremental learning (easy addition and deletion of knowledge), and some provenance for its prediction based on the nearest neighbors w.r.t. the input. However, these models often suffer from weaker empirical performance as compared to deep parametric models.

Increasingly, a middle ground combining the two paradigms and retaining the best of both worlds is becoming popular across various domains, ranging from natural language (Das et al., 2021; Wang et al., 2022; Liu et al., 2022; Izacard et al., 2022), to vision (Liu et al., 2015; 2019; Iscen et al., 2022; Long et al., 2022), to reinforcement learning (Blundell et al., 2016; Pritzel et al., 2017; Ritter et al., 2020), to even protein structure prediction (Cramer, 2021). In such approaches, given a test input, one first retrieves relevant entries from a data index and then processes the retrieved entries along with the test input to make the final predictions using an ML model. This process is visualized in Fig. 1c.

While classical learning setups (cf. Fig. 1a and 1b) have been studied extensively over decades, even basic properties and trade-offs pertaining to retrieval-based models (cf. Fig. 1c), despite their aforementioned remarkable successes, remain highly under-explored. Most of the existing efforts on retrieval-based models solely focus on developing end-to-end *domain-specific* models, without identifying the key dataset properties or structures that are critical in realizing performance gains by such models. Furthermore, at first glance, due to the highly dependent nature of an input and the associated retrieved set, direct application of existing statistical learning techniques does not appear as straightforward. This prompts a natural question:

What is the right theoretical framework to rigorously showcase the value of the retrieved set in ensuring superior performance of modern retrieval-based models?

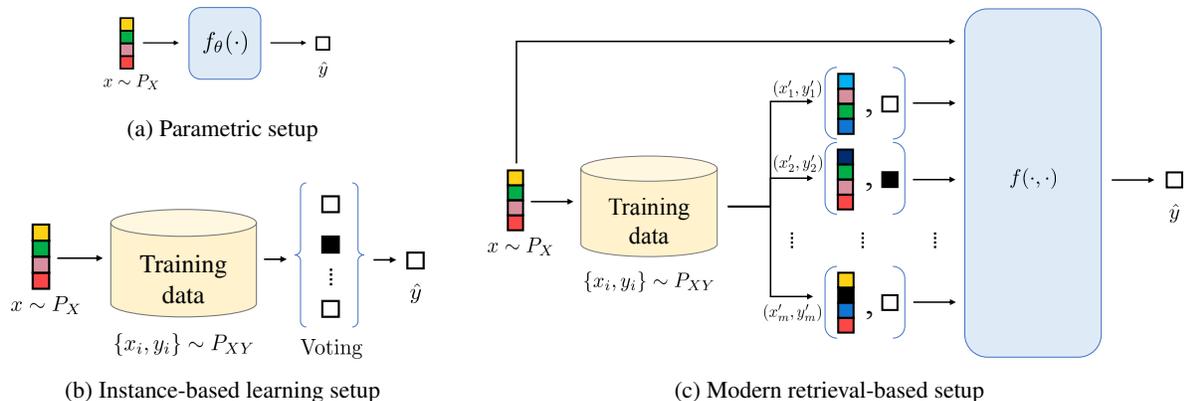


Figure 1. An illustration of a retrieval-based classification model. Given an input instance x , similar to an instance-based model, it retrieves similar (labeled) examples $\mathcal{R}^x = \{(x'_j, y'_j)\}_j$ from training data. Subsequently, it processes input instance along with the retrieved examples (potentially via a nonparametric method) to make the final prediction $\hat{y} = f(x, \mathcal{R}^x)$.

In this paper, we take the first step towards answering this question, while focusing on the classification setting (Sec. 2.1). We begin with the hypothesis that the model might be using the retrieved set to do *local learning* implicitly and then adapt its predictions to the neighborhood of the test point. Multiple recent works (Garg et al., 2022; Akyürek et al., 2022; von Oswald et al., 2022) have studied the feasibility of such a mechanism in widely popular Transformer models. Notably, these works show that a Transformer network can emulate gradient descent to optimize a local learning objective when presented with multiple labeled examples as inputs. Such local learning is potentially beneficial in cases where the underlying task has a local structure, where a much simpler function class suffices to explain the data in a given local neighborhood but overall the data can be complex (formally defined in Sec. 2.2). For example, to solve an issue at hand (a problem instance), it is often faster to search for solutions to similar problems on Stackoverflow and utilize those (i.e., locally learning from the retrieved similar labeled examples) than understanding the whole system (i.e., learning the entire global function).

Inspired by Bottou & Vapnik (1992), we analyze an *explicit* local learning framework: For each test input, 1) we retrieve a few (labeled) training examples located in the vicinity of the test input, 2) train a local model by performing empirical risk minimization (ERM) with only these retrieved examples – *local ERM*; and 3) apply the resulting local model to make prediction on the test input. For the aforementioned retrieval-based local ERM, we derive finite sample generalization bounds that highlight a trade-off between the complexity of the underlying function class and size of the neighborhood where local structure of the data distribution holds in Sec. 3. Under this assumption of local regularity, we show that by using a much simpler function class for the local model, we can achieve a similar loss/error to that of a complex global model (Thm. 3.7). Thus, we show that breaking down the underlying learning task into local sub-tasks enables the

model to employ a *low complexity* parametric component to ensure good global accuracy via a retrieval-based model.

We acknowledge that such local learning cannot be the complete picture behind the effectiveness of retrieval-based models. As noted in Zakai & Ritov (2008), there always exists a model with global component that is more “preferable” to a local-only model. In Sec. 3.4, we extend local ERM to a two-stage setup: First learn a global representation using entire dataset, and then utilize the representation at the test time while solving the local ERM as previously defined. This enables the local learning to benefit from good quality global representations, especially in sparse data regions.

Finally, we move beyond explicit local learning to a setting that resembles more closely the empirically successful systems such as REINA (Wang et al., 2022), WebGPT (Nakano et al., 2021), and AlphaFold (Cramer, 2021): A model that directly learns to predict from the input instance and associated retrieved similar examples end-to-end. Towards this, we take a preliminary step in Sec. 4 by studying a novel formulation of classification over an extended feature space (to account for the retrieved examples) by using kernel methods (Deshmukh et al., 2019).

To summarize, our main contributions include: 1) Setting up a formal framework for classification via retrieval-based models under local structure; 2) Finite sample analysis of explicit local learning framework; 3) Comparison with simple parametric and nonparametric paradigms 4) Extending the analysis to a globally learnt model; and 5) Providing the first rigorous treatment of an end-to-end retrieval-based model to study its generalization by using kernel-based learning.

2. Problem setup

We first provide a brief background on (multiclass) classification along with the necessary notations. Subsequently, we discuss the problem setup considered in this paper, which

deals with designing retrieval-based classification models for the data distributions with local structures.

2.1. Multiclass classification

In this work, we restrict ourselves to (multi-class) classification setting, with access to n training examples $\mathcal{S} = \{(x_i, y_i)\}_{i \in [n]} \subset \mathcal{X} \times \mathcal{Y}$, sampled i.i.d. from the data distribution $D := D_{X,Y}$. Let $\mathbb{P}_D(A) := \mathbb{E}_{(X,Y) \sim D} [\mathbf{1}_{\{A\}}]$ for any random variable A .

Given \mathcal{S} , one is interested in learning a classifier $h : \mathcal{X} \rightarrow \mathcal{Y}$ that minimizes miss-classification error. It is common to define a classifier via a scorer $f : x \mapsto (f_1(x), \dots, f_{|\mathcal{Y}|}(x)) \in \mathbb{R}^{|\mathcal{Y}|}$ that assigns a score to each class in \mathcal{Y} for an instance x . For a scorer f , the corresponding classifier takes the form:

$$h_f(x) = \arg \max_{y \in \mathcal{Y}} f_y(x).$$

Given a set of scorers $\mathcal{F}^{\text{global}} \subseteq \{f : \mathcal{X} \rightarrow \mathbb{R}^{|\mathcal{Y}|}\}$, learning a model implies finding a scorer in $\mathcal{F}^{\text{global}}$ that minimizes the miss-classification error or expected 0/1-loss:

$$f_{0/1}^* = \arg \min_{f \in \mathcal{F}^{\text{global}}} \mathbb{P}_D(h_f(X) \neq Y). \quad (1)$$

One typically employs a surrogate loss (Bartlett et al., 2006) ℓ for the miss-classification error $\mathbb{1}_{\{h_f(X) \neq Y\}}$ and aims to minimize the associated population risk:

$$R_\ell(f) = \mathbb{E}_{(X,Y) \sim D} [\ell(f(X), Y)].$$

Since the underlying data distribution D is only accessible via examples in \mathcal{S} , one learns a good scorer by minimizing the (global) *empirical* risk over the function class $\mathcal{F}^{\text{global}}$ as follows:

$$\hat{f} = \arg \min_{f \in \mathcal{F}^{\text{global}}} \frac{1}{n} \sum_{i \in [n]} \ell(f(x_i), y_i). \quad (2)$$

We denote $\hat{R}_\ell(f) := \frac{1}{n} \sum_{i \in [n]} \ell(f(x_i), y_i)$.

2.2. Classification with local structure

In this work, we assume that the underlying data distribution D has a local structure, where a much simpler (parametric) function class suffices to explain the data in each local neighborhood. Formally, for $x \in \mathcal{X}$ and $r > 0$, we define $\mathcal{B}^{x,r} := \{x' \in \mathcal{X} : \mathfrak{d}(x, x') \leq r\}$, an r -radius ball around x , w.r.t. a metric $\mathfrak{d} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. Let $D^{x,r}$ be the data distribution restricted to $\mathcal{B}^{x,r}$, i.e.,

$$D^{x,r}(A) = D(A) / D(\mathcal{B}^{x,r} \times \mathcal{Y}) \quad A \subseteq \mathcal{B}^{x,r} \times \mathcal{Y}. \quad (3)$$

Further, let us define the local population risk of a function f at a given instance $x \in \mathcal{X}$:

$$R_\ell^x(f) = \mathbb{E}_{(X',Y') \sim D^{x,r}} [\ell(f(X'), Y')].$$

Now, the *local structure condition* of the data distribution ensures that, for each $x \in \mathcal{X}$, there exists a low-complexity function class \mathcal{F}^x , with $|\mathcal{F}^x| \ll |\mathcal{F}^{\text{global}}|$, that approximates

the Bayes optimal (w.r.t. $\mathcal{F}^{\text{global}}$) for the *local classification problem* defined by $D^{x,r}$. That is, for a given $\varepsilon_{\mathcal{X}} > 0$ and $\forall x \in \mathcal{X}$, we have that ¹

$$\min_{f \in \mathcal{F}^x} R_\ell^x(f) \leq \min_{f \in \mathcal{F}^{\text{global}}} R_\ell^x(f) + \varepsilon_{\mathcal{X}}. \quad (4)$$

As an example, if $\mathcal{F}^{\text{global}}$ is linear in \mathbb{R}^d (possibly dense) with bounded norm τ , then \mathcal{F}^x can be a simpler function class such as linear in \mathbb{R}^d with sparsity $k \ll d$ and with bounded norm $\tau_x \leq \tau$.

2.3. Retrieval-based classification model

This work focuses on retrieval-based methods that can leverage the aforementioned local structure of the data distribution. In particular, we focus on two such approaches:

Local empirical risk minimization. Given a (test) instance x , the *local* empirical risk minimization (ERM) approach first retrieves a neighboring set $\mathcal{R}^x = \{(x'_j, y'_j)\} \subseteq \mathcal{S}$. Subsequently, it identifies a (local) scorer \hat{f}^x from a ‘simple’ function class $\mathcal{F}^{\text{loc}} \subset \{f : \mathcal{X} \rightarrow \mathbb{R}^{|\mathcal{Y}|}\}$ as follows:

$$\hat{f}^x = \arg \min_{f \in \mathcal{F}^{\text{loc}}} \frac{1}{|\mathcal{R}^x|} \sum_{(x', y') \in \mathcal{R}^x} \ell(f(x'), y'). \quad (5)$$

By convention if $|\mathcal{R}^x| = 0$, $\hat{f}^x \in \mathcal{F}^{\text{loc}}$ is chosen arbitrarily.

Note that the local ERM approach requires solving a local learning task for each test instance. Such a local learning algorithms was introduced by Bottou & Vapnik (1992). Another point worth mentioning here is that (5) employs the same function class \mathcal{F}^{loc} for each x , whereas the local structure assumption (cf. (4)) allows for an instance dependent function class \mathcal{F}^x . We consider \mathcal{F}^{loc} that approximates $\cup_{x \in \mathcal{X}} \mathcal{F}^x$ closely. In particular, we assume that, for some $\varepsilon_{\text{loc}} > 0$, we have $\forall x \in \mathcal{X}$ that

$$\min_{f \in \mathcal{F}^{\text{loc}}} R_\ell^x(f) \leq \min_{f \in \mathcal{F}^x} R_\ell^x(f) + \varepsilon_{\text{loc}}. \quad (6)$$

Continuing with the example following (4), where \mathcal{F}^x is linear with sparsity $k \ll d$ and bounded norm τ_x , one can take \mathcal{F}^{loc} to be linear with the same sparsity k and bounded norm $\tau' < \sup_{x \in \mathcal{X}} \tau_x$.

Classification with extended feature space. We also consider the setting where the scorer f can implicitly solve the local-ERM using retrieved neighboring labeled instances to make the classification prediction. In other words, the scorer directly maps the augmented input $x \times \mathcal{R}^x \in \mathcal{X} \times (\mathcal{X} \times \mathcal{Y})^*$ to per-class scores. One can learn such a scorer over *extended* feature space $\mathcal{X} \times (\mathcal{X} \times \mathcal{Y})^*$ as follows:

$$\hat{f}^{\text{ex}} = \arg \min_{f \in \mathcal{F}^{\text{ex}}} \hat{R}_\ell^{\text{ex}}(f), \quad (7)$$

where $\hat{R}_\ell^{\text{ex}}(f) := \frac{1}{n} \sum_{i \in [n]} \ell(f(x_i, \mathcal{R}^{x_i}), y_i)$ and a function class of interest over the extended space is denoted

¹As stated, we require the local structure condition to hold for each x . This can be relaxed to hold with high probability with the increased complexity of exposition.

as $\mathcal{F}^{\text{ex}} \subset \{f : \mathcal{X} \times (\mathcal{X} \times \mathcal{Y})^* \rightarrow \mathbb{R}^{|\mathcal{Y}|}\}$. Examples of such a function class include prompting transformers with the retrieved labeled examples. Moreover, it has been recently shown that a transformer can express certain algorithms for optimizing a local learning objective based on the examples from the prompt using gradient descent (Garg et al., 2022; Akyürek et al., 2022; von Oswald et al., 2022).

Our goal is to develop a statistical understanding of these two retrieval-based methods for classification when the underlying data distribution has local structure. We present our theoretical treatment of local ERM and classification with extended feature space in Sec. 3 and 4, respectively.

3. Local empirical risk minimization

In this section, our objective is to characterize the excess risk of local ERM. In particular, we aim to bound

$$\mathbb{E}_{(X,Y) \sim \mathcal{D}}[\ell(\hat{f}^X(X), Y) - \ell(f^*(X), Y)]. \quad (8)$$

Note that \hat{f}^X (cf. (5)) in the above equation is a function of \mathcal{R}^X , and expectation over \mathcal{R}^X is taken implicitly.

3.1. Assumptions

Before presenting an excess risk bound for the local ERM method, we introduce various necessary definitions and assumptions that play a critical role in our analysis.

We define the *margin* of scorer f at a given label $y \in \mathcal{Y}$ as

$$\gamma_f(x, y) = f_y(x) - \max_{y' \neq y} f_{y'}(x). \quad (9)$$

In order to ensure the margin of the scorer f has smooth deviation as x varies, we introduce L -coordinate Lipschitz condition: A scorer f is L -coordinate Lipschitz iff for all $y \in \mathcal{Y}$ and $x, x' \in \mathcal{X}$, we have

$$|f_y(x) - f_y(x')| \leq L \|x - x'\|_2. \quad (10)$$

Following Döring et al. (2018), we define the *weak margin condition* for a scorer f : Given a distribution \mathcal{D} , a scorer f satisfies (α, c) -weak margin condition iff, for all $t \geq 0$,

$$\mathbb{P}_{(X,Y) \sim \mathcal{D}}(|\gamma_f(X, Y)| \leq t) \leq c t^\alpha. \quad (11)$$

One of the key assumptions that we rely on is the existence of an underlying scorer f^{true} that explains the true labels, while ensuring the weak margin condition. Here, we note that the true function f^{true} may neither lie in the function class $\mathcal{F}^{\text{global}}$, nor in \mathcal{F}^{loc} .

Assumption 3.1 (True scorer function). *There exists a scorer f^{true} such that, for all $(x, y) \in \mathcal{X} \times \mathcal{Y}$, f^{true} generates the true label, i.e., $\gamma_{f^{\text{true}}}(x, y) > 0$. Furthermore, we assume f^{true} is L_{true} -coordinate Lipschitz, and satisfies the $(\alpha_{\text{true}}, c_{\text{true}})$ -weak margin condition.*

Furthermore, we restrict ourselves to smooth loss functions that act on the margin of a scorer (cf. (9)).

Assumption 3.2 (Margin-based Lipschitz loss). *For any given example (x, y) and any scorer f , we have $\ell(f(x), y) = \ell(\gamma_f(x, y))$ and ℓ is a decreasing function of the margin. Furthermore, the loss function ℓ is L_ℓ -Lipschitz function, i.e., $|\ell(\gamma) - \ell(\gamma')| \leq L_\ell |\gamma - \gamma'|$, $\forall \gamma \geq \gamma'$.*

Recall that \mathcal{R}^x corresponds to the samples in \mathcal{S} that belong to $\mathcal{B}^{x,r}$; hence, it follows the distribution $\mathcal{D}^{x,r}$. For the rest of the paper, we limit ourselves to $\mathcal{X} \subseteq \mathbb{R}^d$. We can extend this to more general metric spaces with the increased complexity of exposition. Let the density of the distribution of $x \in \mathcal{X} \subseteq \mathbb{R}^d$ be $\rho_{\mathcal{D}}(x)$. A common assumption in the nonparametric estimation literature is the weak density condition (see, e.g., Döring et al., 2018). Moreover, we need to ensure that with high probability the density $\rho_{\mathcal{D}}(x)$ is not too low. We do so following the idea of density level sets from Steinwart (2011). Accordingly, we make the following assumption.

Assumption 3.3 (Data regularity condition).

- (Weak density condition) *There exists constants $c_{\text{wdc}} > 0$, and $\delta_{\text{wdc}} > 0$, such that for all $x \in \mathcal{X}$ and $\rho_{\mathcal{D}}(x)r^d \leq \delta_{\text{wdc}}^d$,*

$$\mathbb{P}_{X' \sim \mathcal{D}}[\text{d}(X', x) \leq r] \geq c_{\text{wdc}}^d \rho_{\mathcal{D}}(x)r^d.$$

- (Density level-set) *There exists a function $f_\rho(\delta)$ with $f_\rho(\delta) \rightarrow 0$ as $\delta \rightarrow 0$, such that for any $\delta > 0$,*

$$\mathbb{P}_{X \sim \mathcal{D}}[\rho_{\mathcal{D}}(X) \leq f_\rho(\delta)] \leq \delta. \quad (12)$$

For example, for d -dimensional multivariate Gaussian with the covariance matrix Σ , we have $f_\rho(\delta) = \Theta(2^{-d/2} |\Sigma|^{-1/2} \delta \ln(1/\delta)^{-d/2})$. This result can be extended to mixture of Gaussian and sub-gaussian random variables (see Appendix B.6 for details).

Assumption 3.4 (Weak+ density condition). *There exists constants $c_{\text{wdc}+} \geq 0$, and $\alpha_{\text{wdc}+} > 0$, such that for all $x \in \mathcal{X}$ and $r \in [0, r_{\text{max}}]$,*

$$\left| \frac{\mathbb{P}_{X' \sim \mathcal{D}}[\text{d}(X', x) \leq r]}{\rho_{\mathcal{D}}(x) \text{vol}_d(r)} - 1 \right| \leq c_{\text{wdc}+} r^{\alpha_{\text{wdc}+}}.$$

The above assumption implies Assumption 3.3.1. We will show that under Assumption 3.4 the local ERM error bounds can be tightened further. For example, in d -dimensional multivariate Gaussian with the covariance matrix Σ , we have $c_{\text{wdc}+} = \frac{d \lambda_{\text{max}}(\Sigma^{-1})}{2(d+2)}$, and $\alpha_{\text{wdc}+} = 2$ for $r_{\text{max}} = \sqrt{(d+2) \lambda_{\text{max}}(\Sigma)}$, where $\lambda_{\text{max}}(\cdot)$ denotes the maximum eigenvalue.

3.2. Excess risk bound for local ERM

We now proceed to our main results on the excess risk bound of local ERM. Recall that, at $x \in \mathcal{X}$, $f^{x,*}$ denotes the minimizer of the population version of the local loss, and f^* the population risk minimizer for the global loss, i.e.,

$$f^{x,*} = \arg \min_{f \in \mathcal{F}^{\text{loc}}} \mathcal{R}_\ell^x(f); \quad f^* = \arg \min_{f \in \mathcal{F}^{\text{global}}} \mathcal{R}_\ell(f). \quad (13)$$

To bound the excess risk defined in Eq. (8), we first obtain the following upper bound on (8).

Lemma 3.5 (Risk decomposition). *The expected excess risk of the local ERM solution \hat{f}^X is bounded as*

$$\begin{aligned}
& \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\ell(\hat{f}^X(X), Y) - \ell(f^*(X), Y) \right] \\
& \leq \underbrace{\mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[R_\ell^X(f^{X,*}) - R_\ell^X(f^*) \right]}_{\text{Local vs Global Population Optimal Risk}} \\
& + \underbrace{\sum_{\mathcal{F} \in \{\mathcal{F}^{\text{global}}, \mathcal{F}^{\text{loc}}\}} \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\sup_{f \in \mathcal{F}} |R_\ell^X(f) - \ell(f(X), Y)| \right]}_{\text{Global and Local: Sample vs Retrieved Set Risk}} \\
& + \underbrace{\mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\sup_{f \in \mathcal{F}^{\text{loc}}} |R_\ell^X(f) - \hat{R}_\ell^X(f)| \right]}_{\text{Generalization of Local ERM}} \\
& + \underbrace{\mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[|R_\ell^X(f^{X,*}) - \hat{R}_\ell^X(f^{X,*})| \right]}_{\text{Central Absolute Moment of } f^{X,*}}.
\end{aligned}$$

We delegate the proof of Lem. 3.5 to Appendix B. Now, as a strategy to obtain desired excess risk bounds, we separately bound the four terms appearing in Lem. 3.5. Note that the first term captures the expected difference between the loss incurred by global population optima $f^* \in \mathcal{F}^{\text{global}}$ and the local population optima $f^{X,*} \in \mathcal{F}^{\text{loc}}$ in a local region around the test instance x . The second term aims to capture the loss for a scorer evaluated at x vs. the expected value of the loss for the scorer at a random instance sampled in the local region of x based on $\mathcal{D}^{x,r}$. The third term corresponds to the standard ‘generalization error’ for the local ERM with respect to the local data distribution $\mathcal{D}^{X,r}$, whereas the fourth term is the empirical variation of the local population optima $f^{X,*}$ around its population mean under $\mathcal{D}^{X,r}$.

Let the coordinate-Lipschitz constants for scorers in \mathcal{F}^{loc} and $\mathcal{F}^{\text{global}}$ be L_{loc} and L_{global} , respectively. We define a function class $\mathcal{G}(X, Y) := \{(x', y') \mapsto \ell(\gamma_f(\cdot, \cdot)) - \ell(\gamma_f(X, Y)) : f \in \mathcal{F}^{\text{loc}}\}$. Here, by subtracting $\ell(f(X), Y)$ from the loss, we center the losses on \mathcal{R}^X for any function $f \in \mathcal{F}^{\text{loc}}$, and obtain a tighter bound by utilizing the local structure of the distribution $\mathcal{D}^{X,r}$. For any $L > 0$, for notational convenience let us define

$$\begin{aligned}
\mathcal{M}_r(L; \ell, f_{\text{true}}, \mathcal{F}) & := \\
& 2L_\ell \left(Lr + (2\|\mathcal{F}\|_\infty - Lr)c_{\text{true}}(2L_{\text{true}}r)^{\alpha_{\text{true}}} \right). \quad (14)
\end{aligned}$$

For any $x \in \mathcal{X}$, the weak density condition provides high probability lower bound on the size of the retrieved set \mathcal{R}^x .

Proposition 3.6. *Under the Assumption 3.3, for any $x \in \mathcal{X}$, radius $r > 0$, and $\delta > 0$,*

$$\mathbb{P}_{\mathcal{D}} \left[|\mathcal{R}^x| < N(r, \delta) \right] \leq \delta, \quad (15)$$

for $N(r, \delta) = n \left(c_{\text{wdc}}^d \min\{f_\rho(\delta/2)r^d, \delta_{\text{wdc}}^d\} - \sqrt{\frac{\log(2/\delta)}{2n}} \right)$.

Now, by controlling different terms appearing in the bound in Lem. 3.5, we obtain the following.

Theorem 3.7 (Excess risk bound). *Let (4) and (6); and Assumptions 3.1, 3.2 and 3.3 hold. For any $\delta > 0$, and $N(r, \delta)$ as defined in Proposition 3.6, the expected excess risk of the local ERM solution \hat{f}^X is bounded as*

$$\begin{aligned}
& \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\ell(\hat{f}^X(X), Y) - \ell(f^*(X), Y) \right] \\
& \leq \underbrace{(\varepsilon_{\mathcal{X}} + \varepsilon_{\text{loc}})}_{\text{Local vs Global Optimal loss (I)}} \\
& + \underbrace{\mathcal{M}_r(L_{\text{loc}}; \ell, f_{\text{true}}, \mathcal{F}^{\text{loc}}) + \mathcal{M}_r(L_{\text{global}}; \ell, f_{\text{true}}, \mathcal{F}^{\text{global}})}_{\text{Global and Local: Sample vs Retrieved Set Risk (II)}} \\
& + \underbrace{\left[\mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\mathfrak{R}_{\mathcal{R}^X}(\mathcal{G}(X, Y)) \mid |\mathcal{R}^X| \geq N(r, \delta) \right] \right.}_{\text{Generalization of Local ERM and Central Absolute Moment of } f^{X,*} \text{ (III)}} \\
& \quad \left. + 5\mathcal{M}_r(L_{\text{loc}}; \ell, f_{\text{true}}, \mathcal{F}^{\text{loc}}) \sqrt{\frac{2 \ln(4/\delta)}{N(r, \delta)}} \right. \\
& \quad \left. + 8\delta L_\ell \|\mathcal{F}^{\text{loc}}\|_\infty \right].
\end{aligned}$$

where $\mathfrak{R}_{\mathcal{R}^X}(\mathcal{G}(X, Y))$ denotes the empirical Rademacher complexity of $\mathcal{G}(X, Y)$. Under Assumption 3.4 and $r \leq r_{\text{max}}$, Sample vs Retrieved Set Risk (II) is $O(c_{\text{wdc}+} r^{\alpha_{\text{wdc}+}})$.

The above result shows a trade-off in *approximation* vs. *generalization* error as retrieval radius r varies.

Approximation error. It comprises two components, defined by (I) and (II) in Thm. 3.7. $\varepsilon_{\mathcal{X}}$ shows the gap in approximating the r -radius neighborhood around X with a simple local function class \mathcal{F}^X which varies with $X \in \mathcal{X}$. Next, ε_{loc} shows the gap in approximating the union of the local function class $\cup_{x \in \mathcal{X}} \mathcal{F}^x$ with a single function class \mathcal{F}^{loc} (possibly with smaller complexity) but while allowing for choosing a different optimizer $f^X \in \mathcal{F}^{\text{loc}}$ for each $X \in \mathcal{X}$. Both $\varepsilon_{\mathcal{X}}$ and ε_{loc} typically increases with r .

The second component of the approximation error (II) corresponds to the difference of risk for the sample X and the retrieved set \mathcal{R}^X for $\mathcal{F}^{\text{global}}$ and \mathcal{F}^{loc} , i.e., $\mathcal{M}_r(L_{\text{global}}; \ell, f_{\text{true}}, \mathcal{F}^{\text{global}})$ and $\mathcal{M}_r(L_{\text{loc}}; \ell, f_{\text{true}}, \mathcal{F}^{\text{loc}})$. Eq. (14) suggests that the terms increase as $O(\text{poly}(r))$. When the data follows multivariate Gaussian then term (II) increases as $O(r^2)$.

Generalization error. It (III) depends on the size of the retrieved set \mathcal{R}^X and the Rademacher complexity of $\mathcal{G}(X, Y)$ which is induced by \mathcal{F}^{loc} . With increasing radius r , the term $N(r, \delta)$ increases. The Rademacher complexity decays with increasing radius, r , typically at the rate of $O(1/\sqrt{N(r, \delta)})$. Thus, under the local ERM setting the total approximation error increases with increasing radius r , given \mathcal{F}^{loc} is fixed. On the contrary, the generalization error decreases with increasing radius r for a fixed \mathcal{F}^{loc} . This suggests a trade-off between the approximation and generalization error as we make a design choice about r . (We empirically validate this in Fig. 3.) Due to centering within the set $\mathcal{G}(X, Y)$ we have

the upper bound on this term as $\mathcal{M}_r(L_{\text{loc}}; \ell, f_{\text{true}}, \mathcal{F}^{\text{loc}})$, which is effective for small r . This does not decay with $|\mathcal{R}^X|$, hence becomes worse with increasing r and complements the above standard.

3.3. Illustrative examples

Assume the $\mathcal{F}^{\text{global}}$ admits q^x -th order derivative in the region $\mathcal{B}(x, r)$. Then a natural choice for $|\mathcal{F}^x|$ is the set of multivariate polynomial functions of degree q^x , namely $\mathcal{P}(q^x)$, for some $q^x \geq 1$. The L_1 approximation error between $\mathcal{F}^x \equiv \mathcal{P}(q^x)$ and $\mathcal{F}^{\text{global}}$ can be quantified using the remainder in Taylor's approximation. This remainder typically grows as $C(\mathcal{F}^{\text{global}}, q^x)r^{(q^x+1)}$ for our choice of radius r for the neighborhood, where $C(\mathcal{F}^{\text{global}}, q^x)$ depends on the function class and the degree. Therefore, we have $\varepsilon_{\mathcal{X}} \leq C(\mathcal{F}^{\text{global}}, q^x)r^{(q^x+1)}$.

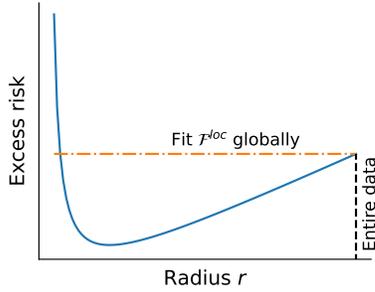


Figure 2. Behavior of excess risk of local ERM

Local linear models. Let us consider this setting where \mathcal{F}^{loc} is the class of linear classifiers in d -dimension. The error in approximating $\mathcal{F}^x = \mathcal{P}(q^x)$ for any $q^x > 1$ with a linear classifier in the $\mathcal{B}(x, r)$ neighborhood for any $x \in \mathcal{X}$ is bounded by $\varepsilon_{\text{loc}} = \Theta(r^2)$. Therefore, the term (I) admits the bound $O(r^2)$. The generalization term varies as $O(1/\sqrt{N(r, \delta)})$. For $r \geq \Omega(n^{-1/2d} \log(n)^{1/2})$ and $\delta = n^{-1/2d}$ then $N(r, \delta) = \Omega(\sqrt{n^{(2d-1)/2d} r^d})$. Combining this we obtain:

$$\begin{aligned} \text{Excess Risk} &\leq \underbrace{O(r^2)}_{\text{(I)}} + \underbrace{O(r^{\min\{\alpha_{\text{true}}, 1\}})}_{\text{(II)}} + \\ &\underbrace{O\left(\frac{d}{n^{(2d-1)/2d} r^{d/2}} + \frac{r^{\min\{\alpha_{\text{true}}, 1\}}}{n^{(2d-1)/4d} r^{d/2}} + \frac{1}{n^{1/2d}}\right)}_{\text{(III)}}. \end{aligned}$$

For $r = n^{-1/2d} \log(n)^{1/2}$ the excess risk bound is $O(n^{-1/2d} \log(n)^{1/2})$, where the bottleneck comes from the term (II), i.e., the sample vs retrieved risk. This is depicted in Fig. 2. Moreover, when the data has multivariate Gaussian distribution we have the term (II) scale as $O(r^2)$, leading to excess risk of $O(n^{-1/d} \log(n)^{1/2})$. However, global ERM with linear classifiers increases the approximation error considerably. In particular, now approximation error becomes a constant $O(\text{diam}(\mathcal{X})^2)$, and dwarfs the generalization that decreases as $O(1/\sqrt{n})$.

Feed-forward classifiers. As another concrete example we study the setting where \mathcal{F}^{loc} is the class of fully connected deep neural networks (FC-DNN). We have $f_y(\cdot)$ to be an L layer feed-forward network with 1-Lipschitz nonlinearities (Bartlett et al., 2017). Let, for layers $l = 1$ to L , the dimension of the weight matrix be $(d_l \times d_{l-1})$ with $d_L = |Y|$. Also, let b_l and s_l be the $\ell_{2,1}$ norm and spectral norm upper bounds for layer l weight matrix, respectively, with $b_l/s_l \leq \kappa$. We define $d_{\text{max}} = \max_{l \in [L]} d_l$ and let $\tilde{B} = \max_{x \in \mathcal{X}} \|x\|_2 \prod_{l=1}^L s_l$.

Approximation Error ε_{loc} : For bounding ε_{loc} in (6), we require L_1 error of \mathcal{F}^{loc} in approximating polynomials of degree $q_{\text{max}} = \max_{x \in \mathcal{X}} q^x$. An FC-DNN that can approximate polynomials with degree at most q_{max} upto L_1 error ε_{loc} has (see, Theorem 9 in Liang & Srikant (2016))²

$$\begin{aligned} \text{depth, } L &= O(q_{\text{max}} + \log(dq_{\text{max}}C'(\mathcal{F}^{\text{global}}, q_{\text{max}})/\varepsilon_{\text{loc}})), \\ \text{width, } d_{\text{max}} &= O(d \log(dq_{\text{max}}C'(\mathcal{F}^{\text{global}}, q_{\text{max}})/\varepsilon_{\text{loc}})). \end{aligned}$$

Here, $C'(\mathcal{F}^{\text{global}}, q^x)$ is a constant independent of r and ε .

Rademacher complexity: We now bound the term $\mathbb{E}_{(X, Y) \sim \mathcal{D}}[\mathfrak{R}_{\mathcal{R}^X}(\mathcal{G}(X, Y)) | |\mathcal{R}^X| > N(r, \delta)]$ for this class. Following (Bartlett et al., 2017), for some universal constant $C'' > 0$ and any $\delta > 0$, we can bound the term as

$$C'' \left(\frac{L \tilde{B} \sqrt{\kappa} \ln(d_{\text{max}}) L^{3/4} \ln(L \tilde{B} \sqrt{n})^{3/2}}{\sqrt{N(r, \delta)}} + 2\delta \tilde{B} \right).$$

We now provide an excess risk bound when \mathcal{F}^{loc} is the class of FC-DNN. Let $r \geq \Omega(n^{-1/2d} \log(n)^{1/2})$ and $\delta = n^{-1/2d}$. Then, $N(r, \delta) = \Omega(\sqrt{n^{(2d-1)/2d} r^d})$. Now, by setting $\varepsilon_{\text{loc}} = r^{(q_{\text{max}}+1)}$, it follows from Thm. 3.7 that

$$\begin{aligned} \text{Excess Risk} &\leq \underbrace{O(r^{(q_{\text{max}}+1)})}_{\text{(I)}} + \underbrace{O(r^{\min\{\alpha_{\text{true}}, 1\}})}_{\text{(II)}} + \\ &\underbrace{O\left(\frac{q_{\text{max}}^{3/4} \ln(dq_{\text{max}}/r)^{3/4} \ln(n)^{3/2}}{n^{(2d-1)/2d} r^{d/2}} + \frac{r^{\min\{\alpha_{\text{true}}, 1\}}}{n^{(2d-1)/4d} r^{d/2}} + \frac{1}{n^{1/2d}}\right)}_{\text{(III)}}. \end{aligned}$$

With $r = n^{-1/2d} \log(n)^{1/2}$, the excess risk is bounded as $O(n^{-1/2d} \log(n)^{1/2})$. Again (II) is the bottleneck. This bottleneck can be improved for multivariate Gaussian distribution with excess risk $O(n^{-1/d} \log(n)^{1/2})$.

Note that, it's also worth comparing local-ERM with conventional (non-local) ERM. Under the local structure condition (Sec. 2.2), one would utilize a simple \mathcal{F}^{loc} for local-ERM. This would correspond to the Rademacher complexity term in Thm. 3.7 being small. In contrast, the generalization bound for the traditional (non-local) ERM approach would depend on the Rademacher complexity of a function class $\mathcal{F}^{\text{global}}$ that can achieve a low approximation error on the *entire domain*. Such a function class (even under the local structure assumption) would be much more complex than

²Although width is not explicitly mentioned in (Liang & Srikant, 2016), it can be inferred from the constructions.

\mathcal{F}^{loc} , resulting in a large Rademacher complexity. For the right design choice of r , and \mathcal{F}^{loc} , the approximation error increase of local-ERM can be offset by large generalization error of $\mathcal{F}^{\text{global}}$. As a consequence, local ERM with simple function class \mathcal{F}^{loc} can outperform (non-local) ERM with a complex class $\mathcal{F}^{\text{global}}$.

3.4. Endowing local ERM with global representations

Note that the local ERM method takes a somewhat myopic view and does not aim to learn a global hypothesis that (partially or entirely) explains the entire data distribution. Such an approach may potentially result in poor performance in those regions of input domains that are not well represented in the training set. Here, we explore a two-stage approach leveraging the global pattern present in the training data to address this apparent shortcoming of local ERM.

Given training data \mathcal{S} and a simple function class $\mathcal{G}^{\text{loc}} : \mathbb{R}^d \rightarrow \mathbb{R}^{|\mathcal{Y}|}$, the first stage involves learning a d -dimensional feature map $\Phi_{\mathcal{S}} : \mathcal{X} \rightarrow \mathbb{R}^d$ that simultaneously ensures good representation for the entire data distribution (Radford et al., 2021; Grill et al., 2020; Cer et al., 2018; Reimers & Gurevych, 2019). Subsequently, given a test instance x and its retrieved neighboring points $\mathcal{R}^x = \{(x'_j, y'_j)\} \subseteq \mathcal{S}$, one employs local ERM with the function class:

$$\mathcal{F}_{\Phi_{\mathcal{S}}} = \{x \mapsto g \circ \Phi_{\mathcal{S}}(x) : g \in \mathcal{G}^{\text{loc}}\}. \quad (16)$$

At this point, it is tempting to invoke the proof strategy outlined following Lem. 3.5, with \mathcal{F}^{loc} replaced with $\mathcal{F}_{\Phi_{\mathcal{S}}}$ to characterize the performance of the aforementioned two-stage method. Note that one can indeed bound the first two terms appearing in Lem. 3.5 for the two-stage method as well. However, bounding the third term that corresponds to generalization gap for local ERM becomes challenging as $\mathcal{F}_{\Phi_{\mathcal{S}}}$ depends on \mathcal{S} via the global representation $\Phi_{\mathcal{S}}$ learned in the first stage. Interestingly, Foster et al. (2019) explored a general framework to address such dependence for standard (non retrieval-based) learning. In fact, as an instantiation of their general framework, Foster et al. (2019, Sec. 5.4) consider the ERM in feature space defined by a representation. We employ their techniques to obtain the following result on the generalization gap for local ERM with $\mathcal{F}_{\Phi_{\mathcal{S}}}$.

Proposition 3.8. *Assuming the representation learned during the first stage is Δ -sensitive, i.e., for \mathcal{S} and \mathcal{S}' that differ in a single example, we have $\|\Phi_{\mathcal{S}}(x) - \Phi_{\mathcal{S}'}(x)\| \leq \Delta \forall x \in \mathcal{X}$. Furthermore, we assume that each $g \in \mathcal{G}^{\text{loc}}$ (cf. 16) is L -Lipschitz, the loss $\ell : \mathbb{R}^{|\mathcal{Y}|} \times |\mathcal{Y}| \rightarrow \mathbb{R}$ is $L_{\ell,1}$ -Lipschitz w.r.t. $\|\cdot\|_{\infty}$ -norm in the first argument, and ℓ is bounded by M_{ℓ} . Then, following holds with probability at least $1 - \delta$:*

$$\begin{aligned} & \sup_{f \in \mathcal{F}_{\Phi_{\mathcal{S}}}} \left| \mathbb{E}_{(X', Y') \sim \mathcal{D}^{x,r}} [\ell(f(X'), Y')] - \hat{R}_{\ell}^x(f) \right| \\ & \leq (M_{\ell} + 2\Delta L L_{\ell,1} |\mathcal{R}^x|) \sqrt{\log(1/\delta)/2|\mathcal{R}^x|} + \quad (17) \\ & \quad \mathbb{E}_{\mathcal{R}^x \sim \mathcal{D}^{x,r}} \left[\sup_{f \in \mathcal{F}_{\Phi_{\mathcal{S}}}} |R_{\ell}(f) - \hat{R}_{\ell}^x(f)| \right]. \end{aligned}$$

Furthermore,

$$\mathbb{E}_{\mathcal{R}^x \sim \mathcal{D}^{x,r}} \left[\sup_{f \in \mathcal{F}_{\Phi_{\mathcal{S}}}} |R_{\ell}(f) - \hat{R}_{\ell}^x(f)| \right] \leq 2\mathfrak{N}^{\circ}(\ell \circ \mathcal{F}_{\Phi_{\mathcal{S}}}), \quad (18)$$

where $\ell \circ \mathcal{F}_{\Phi_{\mathcal{S}}} = \{(x, y) \mapsto \ell(f(x), y) : f \in \mathcal{F}_{\Phi_{\mathcal{S}}}\}$ and \mathfrak{N}° denotes the Rademacher complexity of data dependent hypothesis sets (Foster et al., 2019).

We defer the proof of Prop. 3.8 and necessary background on Foster et al. (2019) to Appendix D.

As a potential advantage of utilizing a global representation with local ERM, one can realize high-performance local learning with an even simpler function class. For example, it's a common approach to only train a linear classifier on learned representations. Furthermore, a high-quality global representation can ensure good performance for those local regions that are not well represented in the training set. We leave a formal treatment of these topics for a longer version of this manuscript.

4. Classification in extended feature space

Next, we focus on a family of retrieval-based methods that directly learn a scorer to map an input instance and its neighboring labeled instance to a score vector (cf. (7)). In fact, as discussed in Sec. 1, many successful modern instances of retrieval-based models such as REINA (Wang et al., 2022) and KATE (Liu et al., 2022) belong to this family. In this section, we provide the first rigorous treatment (to the best of our knowledge) for such models.

As introduced in Sec. 2.3, our objective is to learn a function $f : \mathcal{X} \times (\mathcal{X} \times \mathcal{Y})^* \rightarrow \mathbb{R}^{|\mathcal{Y}|}$. For a given instance x , such a function can leverage its neighboring set $\mathcal{R}^x \in (\mathcal{X} \times \mathcal{Y})^*$ to improve the prediction on x . In this work, we restrict ourselves to a sub-family of such retrieval-based methods that first map $\mathcal{R}^x \sim \mathcal{D}^{x,r}$ to $\hat{\mathcal{D}}^{x,r}$ — an empirical estimate of the local distribution $\mathcal{D}^{x,r}$, which is subsequently utilized to make a prediction for x . In particular, the scorers of interest are of the form: $(x, \mathcal{R}^x) \mapsto f(x, \hat{\mathcal{D}}^{x,r})$, with

$$f(x, \hat{\mathcal{D}}^{x,r}) = (f_1(x, \hat{\mathcal{D}}^{x,r}), \dots, f_{|\mathcal{Y}|}(x, \hat{\mathcal{D}}^{x,r})) \in \mathbb{R}^{|\mathcal{Y}|}.$$

Here, $f_y(x, \hat{\mathcal{D}}^{x,r})$ denotes the score assigned to the y -th class. Thus, assuming that $\Delta_{\mathcal{X} \times \mathcal{Y}}$ denotes the set of distribution over $\mathcal{X} \times \mathcal{Y}$, we restrict to a suitable function class in $\{f : \mathcal{X} \times \Delta_{\mathcal{X} \times \mathcal{Y}} \rightarrow \mathbb{R}^{|\mathcal{Y}|}\}$. Note that, given a surrogate loss $\ell : \mathbb{R}^{|\mathcal{Y}|} \times \mathcal{Y} \rightarrow \mathbb{R}$ and scorer f , the empirical risk $\hat{R}_{\ell}^{\text{ex}}(f)$ and population risk $R_{\ell}^{\text{ex}}(f)$ take the following form:

$$\hat{R}_{\ell}^{\text{ex}}(f) = \frac{1}{n} \sum_{i \in [n]} \ell(x_i, \hat{\mathcal{D}}^{x_i, r})$$

and

$$R_{\ell}^{\text{ex}}(f) = \mathbb{E}_{(X, Y) \sim \mathcal{D}} [\ell(f(X, \mathcal{D}^{X,r}), Y)].$$

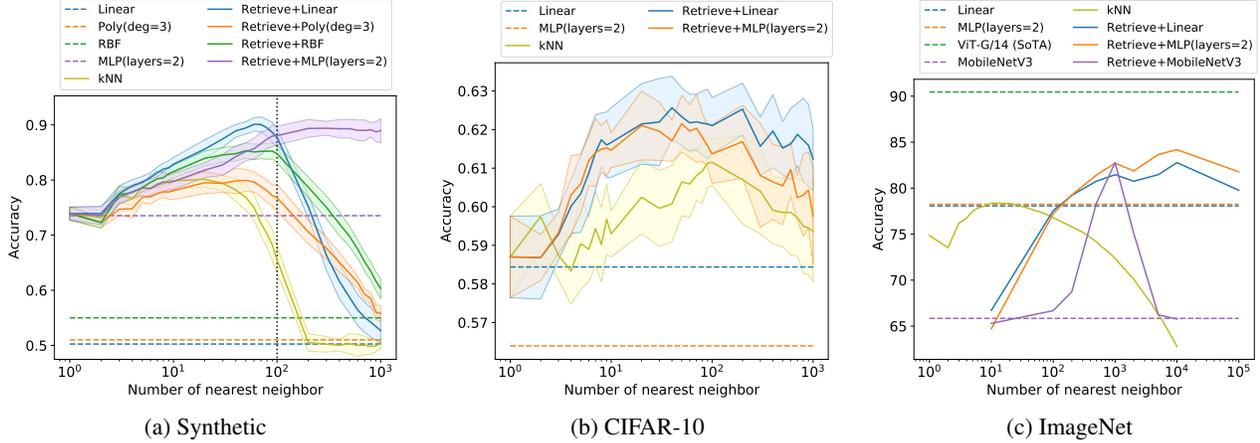


Figure 3. Performance of local ERM with size of retrieved set across models of different complexity.

Note that the general framework for learning in the extended feature space $\tilde{\mathcal{X}} := \mathcal{X} \times \Delta_{\mathcal{X} \times \mathcal{Y}}$ provides a very rich class of functions. In this paper, we focus on a specific form of learning methods in the extended feature space by using the kernel methods. The method as well as its analysis is obtained by adapting the work on utilizing kernel methods for domain generalization (Blanchard et al., 2011; Deshmukh et al., 2019).

In particular, we study generalization of a kernel-based classifier over $\tilde{\mathcal{X}}$ learnt via regularized ERM. Due to space constraints, we present an informal version of our result below. See Appendix E for the precise statement (cf. Thm. E.4), necessary background, and detailed proof.

Theorem 4.1 (Informal). *Let $0 \leq \delta \leq 1$ and $N(r, \delta)$ be as defined in (15). Then, under appropriate assumptions, with probability at least $1 - \delta$, we have*

$$\sup_{f \in \mathcal{F}} |\widehat{R}_\ell^{\text{ex}}(f) - R_\ell^{\text{ex}}(f)| \lesssim C_1 n^{-\frac{1}{2}} \left(1 + \log^{\frac{3}{2}} \sqrt{2n|y|}\right) + C_2 \sqrt{\frac{\log(\frac{2}{\delta})}{N(r, \frac{\delta}{n})}} + C_3 \sqrt{\frac{\log(\frac{1}{\delta})}{n}},$$

where \mathcal{F} is extended feature kernel function class; and $\widehat{R}_\ell^{\text{ex}}(f)$ and $R_\ell^{\text{ex}}(f)$ are empirical and population risks.

Interestingly, the bound in Thm. 4.1 implies that the size of the retrieved set \mathcal{R}^x (as captured by $N(r, \frac{\delta}{n})$) has to scale at least logarithmically in the size of the training set n to ensure convergence.

5. Experiments

There have been numerous successful practical applications of retrieval-based models in the literature (e.g., Wang et al., 2022; Das et al., 2021). Here, we present a brief empirical study for such models in order to corroborate the benefits predicted by our theoretical results. We also present preliminary experiments to empirically verify the kernel based

extended feature space-based approach in Appendix E.3.

Task and dataset. We perform experiments on both synthetic and real datasets, as summarized below. Further details are relegated to Appendix F.

(i) *Synthetic.* We consider a task of binary classification on a Gaussian mixture. Each mixture component is endowed with its local linear decision boundary. We randomly generate a train set of size $n = 10000$ in a 10-dimensional space. We use Euclidean distance for retrieval and perform a 10-fold cross-validation.

(ii) *CIFAR-10.* Next, we consider a task of binary classification on a *real data* for object detection. In particular, we consider a subset of CIFAR-10 dataset where we only restrict to images from "Cat" and "Dog" classes. We randomly partition the data into a train set of size $n = 10000$ points and remaining 2000 points for test. We use Euclidean distance for retrieval and do a 10-fold cross-validation.

(iii) *ImageNet.* Finally, we consider 1000-way classification task on ImageNet dataset. We use the standard train-test split with $n = 1281167$ training and 50000 test examples. Following standard practice in literature, we use unsupervised but globally learned features from ALIGN (Jia et al., 2021) to do image retrieval. This also showcases benefits of endowing local ERM with global representation (Sec. 3.4). Given large computational cost, we could only run each experiment once in this setting.

Methods. On all datasets, as baseline, we consider simple linear classifier and multi-layer perceptron (MLP) of two layers. For retrieval-based models, we consider each of the above methods as the local model to fit on retrieved data points via local ERM framework (Sec. 3). For synthetic datasets, we also considered support vector machines with polynomial kernel (of degree 3) and with radial basis function (RBF) kernel, both for baseline and local ERM. For ImageNet, we additionally consider the state-of-the-art (SoTA) single model published for this task, which is

from the most recent CVPR 2022 (Zhai et al., 2022), as a baseline. In addition, for ImageNet, we also consider the pretrain-finetune version of local ERM, where using the retrieved set we fine-tune a MobileNetV3 (Howard et al., 2019) model that has been pretrained on entire ImageNet.

Observations. In Fig. 3, we observe the tradeoff of varying the size of the retrieved set (as dictated by the neighborhood radius) on the performance of retrieval-based methods across all settings. We see that when the number of retrieved samples is small, local ERM has lower accuracy, this is due to large generalization error. When the size of the retrieved sample space is high, local ERM fails to minimize the loss effectively due to the lack of model capacity. We see that this effect being more pronounced for simpler function classes such as linear classifier as compared to MLP. In Fig. 3c, we see that, via local ERM with a small MobileNet-V3 model, we are able to achieve the top-1 accuracy of 82.78 whereas a regularly trained MobileNet-V3 model achieves the top-1 accuracy of only 65.80. Also the result is very competitive with SoTA of 90.45 with a *much larger model*. Thus, our empirical evaluation demonstrates the utility of retrieval-based models via simple local ERM framework. In particular, it allows small sized models to attain very high performance.

6. Related work and discussion

Local polynomial regression. Perhaps the most similar to our setup is the rich set of work on local polynomial regression, which has been around for a long time since the pioneering works of Stone (1977; 1980). This line of work aims to fit a low-degree polynomial at each point in the data set based on a subset of data points. Such approaches gained a lot of attention as parametric regression was not adequate in various practical applications of the time. The performance of this approach critically depends on subset selected to locally fit the data. Towards this, various selection approaches have been considered: fixed bandwidth (Katzkovnik & Kheisin, 1979), nearest neighbors (Cleveland, 1979), kernel weighted (Ruppert & Wand, 1994), and adaptive methods (Ruppert et al., 1995). All these work only analyze under mean squared error loss and do not handle classification nor provide finite sample generalization bounds, which we obtain in this work.

Multi-task and meta learning. At a surface level, our setup might resemble multi-task and meta learning frameworks. In multi-task learning, we are given the examples from T tasks/distributions and the objective is to ensure good classification performance on all the tasks. In meta-learning, the setting is made harder by requiring good performance on a new target task. As a common approach in these settings, we learn a shared representation across the tasks and then learn a simple task-specific mapping on top of these learned shared features (Vilalta & Drissi, 2002, interalia). Theoretical investigations is quite limited: a few works study

upper-bounds of generalization error in multi-task environments (Ben-David & Borbely, 2008; Ben-David et al., 2010; Pentina & Lampert, 2014; Amit & Meir, 2017), and even fewer in case of meta-learning (Balcan et al., 2019; Khodak et al., 2019; Du et al., 2020; Tripuraneni et al., 2021). However, most of these works assume linear or other simple class, whereas we consider general function class using kernel methods. It is not clear if the aforementioned representation based approach can apply to our setting because: each tasks have little overlap, very large number of tasks, and most importantly a priori an example belongs is not assigned to a task. Interestingly, in this work, we show that retrieval-based approach alleviate the needs to identify the task-membership. Here, we would like to highlight a contemporary work (Li et al., 2023) that studies in-context learning by Transformer models in a multi-task/meta-learning setting. In particular, this work relies on the notion of algorithm stability (Bousquet & Elisseeff, 2002) and presents generalization bounds for Transformers as in-context learners.

7. Conclusion and future direction

In this work, we initiate the development of a theoretical framework to study the statistical properties of retrieval-based modern machine learning models. Our treatment of an explicit local learning paradigm, namely local-ERM, establishes an approximation vs. generalization error trade-off. This highlights the advantage realized by access to a retrieved set during classification as it enables good performance with much simpler (local) function classes. As for the retrieval-based models that leverage a retrieved set without explicitly performing local learning, we present a systematic study by considering a kernel-based classifier over extended feature space. Studying end-to-end retrieval-based models beyond kernel-based classification is a natural and fruitful direction for future work. It’s also worth exploring if existing retrieval-based end-to-end models inherently perform *implicit* local learning via architectures such as Transformers.

References

- Akyürek, E., Schuurmans, D., Andreas, J., Ma, T., and Zhou, D. What learning algorithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*, 2022.
- Amit, R. and Meir, R. Meta-learning by adjusting priors based on extended PAC-bayes theory. *arXiv preprint arXiv:1711.01244*, 2017.
- Bai, Y. and Lee, J. D. Beyond linearization: On quadratic and higher-order approximation of wide neural networks. *arXiv preprint arXiv:1910.01619*, 2019.
- Balcan, M.-F., Khodak, M., and Talwalkar, A. Provable guarantees for gradient-based meta-learning. In *Internat-*

- tional Conference on Machine Learning*, pp. 424–433. PMLR, 2019.
- Bartlett, P. L., Jordan, M. I., and McAuliffe, J. D. Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138–156, 2006.
- Bartlett, P. L., Foster, D. J., and Telgarsky, M. J. Spectrally-normalized margin bounds for neural networks. *Advances in neural information processing systems*, 30, 2017.
- Ben-David, S. and Borbely, R. S. A notion of task relatedness yielding provable multiple-task learning guarantees. *Machine learning*, 73(3):273–287, 2008.
- Ben-David, S., Blitzer, J., Crammer, K., Kulesza, A., Pereira, F., and Vaughan, J. W. A theory of learning from different domains. *Machine learning*, 79(1-2):151–175, 2010.
- Blanchard, G., Lee, G., and Scott, C. Generalizing from several related classification tasks to a new unlabeled sample. In Shawe-Taylor, J., Zemel, R., Bartlett, P., Pereira, F., and Weinberger, K. (eds.), *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc., 2011.
- Blundell, C., Uria, B., Pritzel, A., Li, Y., Ruderman, A., Leibo, J. Z., Rae, J., Wierstra, D., and Hassabis, D. Model-free episodic control. *arXiv preprint arXiv:1606.04460*, 2016.
- Bottou, L. and Vapnik, V. Local Learning Algorithms. *Neural Computation*, 4(6):888–900, 11 1992. ISSN 0899-7667. doi: 10.1162/neco.1992.4.6.888.
- Bousquet, O. and Elisseeff, A. Stability and generalization. *Journal of machine learning research*, 2(Mar):499–526, 2002.
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language models are few-shot learners, 2020.
- Cer, D., Yang, Y., Kong, S.-y., Hua, N., Limtiaco, N., John, R. S., Constant, N., Guajardo-Cespedes, M., Yuan, S., Tar, C., et al. Universal sentence encoder. *arXiv preprint arXiv:1803.11175*, 2018.
- Chen, M., Bai, Y., Lee, J. D., Zhao, T., Wang, H., Xiong, C., and Socher, R. Towards understanding hierarchical learning: Benefits of neural representations. *Advances in Neural Information Processing Systems*, 33:22134–22145, 2020.
- Cleveland, W. S. Robust locally weighted regression and smoothing scatterplots. *Journal of the American statistical association*, 74(368):829–836, 1979.
- Cramer, P. Alphafold2 and the future of structural biology. *Nature Structural & Molecular Biology*, 28(9):704–705, 2021.
- Das, R., Zaheer, M., Thai, D., Godbole, A., Perez, E., Lee, J. Y., Tan, L., Polymenakos, L., and McCallum, A. Case-based reasoning for natural language queries over knowledge bases. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 9594–9611, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.755.
- Deshmukh, A. A., Lei, Y., Sharma, S., Dogan, U., Cutler, J. W., and Scott, C. A generalization error bound for multi-class domain generalization, 2019.
- Döring, M., Györfi, L., and Walk, H. Rate of convergence of k -nearest-neighbor classification rule. *Journal of Machine Learning Research*, 18(227):1–16, 2018.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houshly, N. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- Du, S. S., Hu, W., Kakade, S. M., Lee, J. D., and Lei, Q. Few-shot learning via learning the representation, provably. *arXiv preprint arXiv:2002.09434*, 2020.
- Fix, E. and Hodges, J. L. Discriminatory analysis. non-parametric discrimination: Consistency properties. *International Statistical Review/Revue Internationale de Statistique*, 57(3):238–247, 1989.
- Foster, D. J., Greenberg, S., Kale, S., Luo, H., Mohri, M., and Sridharan, K. Hypothesis set stability and generalization. In Wallach, H., Larochelle, H., Beygelzimer, A., d’Alché Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- Garg, S., Tsipras, D., Liang, P., and Valiant, G. What can transformers learn in-context? a case study of simple function classes. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=f1nZJ2eOet>.
- Grill, J.-B., Strub, F., Althé, F., Tallec, C., Richemond, P., Buchatskaya, E., Doersch, C., Avila Pires, B., Guo, Z., Gheshlaghi Azar, M., et al. Bootstrap your own latent: a new approach to self-supervised learning. *Advances*

- in *Neural Information Processing Systems*, 33:21271–21284, 2020.
- Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E., Cai, T., Rutherford, E., Casas, D. d. L., Hendricks, L. A., Welbl, J., Clark, A., et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- Howard, A., Sandler, M., Chu, G., Chen, L.-C., Chen, B., Tan, M., Wang, W., Zhu, Y., Pang, R., Vasudevan, V., Le, Q. V., and Adam, H. Searching for mobilenetv3. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.
- Iscen, A., Fathi, A., Schmid, C., Caron, M., and Bird, T. A memory transformer network for incremental learning. *arXiv preprint*, 2022.
- Izacard, G., Lewis, P., Lomeli, M., Hosseini, L., Petroni, F., Schick, T., Dwivedi-Yu, J., Joulin, A., Riedel, S., and Grave, E. Few-shot learning with retrieval augmented language models. *arXiv preprint arXiv:2208.03299*, 2022.
- Jia, C., Yang, Y., Xia, Y., Chen, Y.-T., Parekh, Z., Pham, H., Le, Q., Sung, Y.-H., Li, Z., and Duerig, T. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*, pp. 4904–4916. PMLR, 2021.
- Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., and Amodei, D. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- Katkovnik, V. Y. and Kheisin, V. Dynamic stochastic approximation of polynomials drifts. *Avtomatika i Telemekhanika*, pp. 89–98, 1979.
- Khodak, M., Balcan, M.-F. F., and Talwalkar, A. S. Adaptive gradient-based meta-learning methods. *Advances in Neural Information Processing Systems*, 32, 2019.
- Lei, Y., Dogan, Ü., Zhou, D.-X., and Kloft, M. Data-dependent generalization bounds for multi-class classification. *IEEE Transactions on Information Theory*, 65(5): 2995–3021, 2019.
- Li, Y., Swersky, K., and Zemel, R. Generative moment matching networks. In *International conference on machine learning*, pp. 1718–1727. PMLR, 2015.
- Li, Y., Ildiz, M. E., Papailiopoulos, D., and Oymak, S. Transformers as algorithms: Generalization and stability in in-context learning. *arXiv preprint arXiv:2301.07067*, 2023.
- Liang, S. and Srikant, R. Why deep neural networks for function approximation? *arXiv preprint arXiv:1610.04161*, 2016.
- Liu, J., Shen, D., Zhang, Y., Dolan, B., Carin, L., and Chen, W. What makes good in-context examples for GPT-3? In *Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures*, pp. 100–114, Dublin, Ireland and Online, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.deelio-1.10. URL <https://aclanthology.org/2022.deelio-1.10>.
- Liu, S., Liang, X., Liu, L., Shen, X., Yang, J., Xu, C., Lin, L., Cao, X., and Yan, S. Matching-cnn meets knn: Quasi-parametric human parsing. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1419–1427, 2015.
- Liu, Z., Miao, Z., Zhan, X., Wang, J., Gong, B., and Yu, S. X. Large-scale long-tailed recognition in an open world. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2537–2546, 2019.
- Long, A., Yin, W., Ajanthan, T., Nguyen, V., Purkait, P., Garg, R., Blair, A., Shen, C., and van den Hengel, A. Retrieval augmented classification for long-tail visual recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6959–6969, 2022.
- Muandet, K., Fukumizu, K., Sriperumbudur, B., Schölkopf, B., et al. Kernel mean embedding of distributions: A review and beyond. *Foundations and Trends® in Machine Learning*, 10(1-2):1–141, 2017.
- Nakano, R., Hilton, J., Balaji, S., Wu, J., Ouyang, L., Kim, C., Hesse, C., Jain, S., Kosaraju, V., Saunders, W., et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- Pentina, A. and Lampert, C. A PAC-bayesian bound for life-long learning. In *International Conference on Machine Learning*, pp. 991–999, 2014.
- Pritzel, A., Uria, B., Srinivasan, S., Badia, A. P., Vinyals, O., Hassabis, D., Wierstra, D., and Blundell, C. Neural episodic control. In *International Conference on Machine Learning*, pp. 2827–2836. PMLR, 2017.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pp. 8748–8763. PMLR, 2021.
- Reimers, N. and Gurevych, I. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*, 2019.

- Ritter, S., Faulkner, R., Sartran, L., Santoro, A., Botvinick, M., and Raposo, D. Rapid task-solving in novel environments. *arXiv preprint arXiv:2006.03662*, 2020.
- Ruppert, D. and Wand, M. P. Multivariate locally weighted least squares regression. *The annals of statistics*, pp. 1346–1370, 1994.
- Ruppert, D., Sheather, S. J., and Wand, M. P. An effective bandwidth selector for local least squares regression. *Journal of the American Statistical Association*, 90(432): 1257–1270, 1995.
- Samarin, M., Roth, V., and Belius, D. On the empirical neural tangent kernel of standard finite-width convolutional neural network architectures. *arXiv preprint arXiv:2006.13645*, 2020.
- Shalev-Shwartz, S. and Ben-David, S. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- Smola, A., Gretton, A., Song, L., and Schölkopf, B. A hilbert space embedding for distributions. In Hutter, M., Servedio, R. A., and Takimoto, E. (eds.), *Algorithmic Learning Theory*, pp. 13–31, Berlin, Heidelberg, 2007. Springer Berlin Heidelberg. ISBN 978-3-540-75225-7.
- Steinwart, I. Adaptive density level set clustering. In *Proceedings of the 24th Annual Conference on Learning Theory*, pp. 703–738. JMLR Workshop and Conference Proceedings, 2011.
- Steinwart, I. and Christmann, A. *Support Vector Machines*. Springer Publishing Company, Incorporated, 1st edition, 2008. ISBN 0387772413.
- Stone, C. J. Consistent nonparametric regression. *The annals of statistics*, pp. 595–620, 1977.
- Stone, C. J. Optimal rates of convergence for nonparametric estimators. *The annals of Statistics*, pp. 1348–1360, 1980.
- Tripuraneni, N., Jin, C., and Jordan, M. Provable meta-learning of linear representations. In *International Conference on Machine Learning*, pp. 10434–10443. PMLR, 2021.
- Vilalta, R. and Drissi, Y. A perspective view and survey of meta-learning. *Artificial intelligence review*, 18(2):77–95, 2002.
- von Oswald, J., Niklasson, E., Randazzo, E., Sacramento, J., Mordvintsev, A., Zhmoginov, A., and Vladymyrov, M. Transformers learn in-context by gradient descent. *arXiv preprint arXiv:2212.07677*, 2022.
- Wang, S., Xu, Y., Fang, Y., Liu, Y., Sun, S., Xu, R., Zhu, C., and Zeng, M. Training data is more valuable than you think: A simple and effective method by retrieving from training data. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3170–3179, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.226. URL <https://aclanthology.org/2022.acl-long.226>.
- Yarotsky, D. Error bounds for approximations with deep relu networks. *Neural Networks*, 94:103–114, 2017.
- Zakai, A. and Ritov, Y. How local should a learning method be?. In *COLT*, pp. 205–216. Citeseer, 2008.
- Zhai, X., Kolesnikov, A., Houlsby, N., and Beyer, L. Scaling vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12104–12113, 2022.
- Zhang, T. Covering number bounds of certain regularized linear function classes. *Journal of Machine Learning Research*, 2(Mar):527–550, 2002.
- Zhang, T. Statistical analysis of some multi-category large margin classification methods. *Journal of Machine Learning Research*, 5(Oct):1225–1251, 2004.

A. Preliminaries

Definition A.1 (Rademacher complexity). *Given a sample $\mathcal{S} = \{z_i = (x_i, y_i)\}_{i \in [n]} \subset \mathcal{Z}$ and a real-valued function class $\mathcal{F} : \mathcal{Z} \rightarrow \mathbb{R}$, the empirical Rademacher complexity of \mathcal{F} with respect to \mathcal{S} is defined as*

$$\mathfrak{R}_{\mathcal{S}}(\mathcal{F}) = \frac{1}{n} \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \sum_{i=1}^n \sigma_i f(z_i) \right], \quad (19)$$

where $\sigma = \{\sigma_i\}_{i \in [n]}$ is a collection of n i.i.d. Bernoulli random variables. For $n \in \mathbb{N}$, the Rademacher complexity $\bar{\mathfrak{R}}_n(\mathcal{F})$ and worst case Rademacher complexity $\mathfrak{R}_n(\mathcal{F})$ are defined as follows.

$$\bar{\mathfrak{R}}_n(\mathcal{F}) = \mathbb{E}_{\mathcal{S} \sim \mathcal{D}^n} [\mathfrak{R}_{\mathcal{S}}(\mathcal{F})], \quad \text{and} \quad \mathfrak{R}_n(\mathcal{F}) = \sup_{\mathcal{S} \sim \mathcal{Z}^n} \mathfrak{R}_{\mathcal{S}}(\mathcal{F}). \quad (20)$$

Definition A.2 (Covering Number). *Let $\epsilon > 0$ and $\|\cdot\|$ be a norm defined over \mathbb{R}^n . Given a function class $\mathcal{F} : \mathcal{Z} \rightarrow \mathbb{R}$ and a collection of points $\mathcal{S} = \{z_i\}_{i \in [n]} \subset \mathcal{Z}$, we call a set of points $\{u_j\}_{j \in [m]} \subset \mathbb{R}^n$ an $(\epsilon, \|\cdot\|)$ -cover of \mathcal{F} with respect to \mathcal{S} , if we have*

$$\sup_{f \in \mathcal{F}} \min_{j \in [m]} \|f(\mathcal{S}) - u_j\| \leq \epsilon, \quad (21)$$

where $f(\mathcal{S}) = (f(z_1), \dots, f(z_n)) \in \mathbb{R}^n$. The $\|\cdot\|$ -covering number $\mathcal{N}_{\|\cdot\|}(\epsilon, \mathcal{F}, \mathcal{S})$ denotes the cardinality of the minimal $(\epsilon, \|\cdot\|)$ -cover of \mathcal{F} with respect to \mathcal{S} . In particular, if $\|\cdot\|$ is an normalized- ℓ_p norm ($\|v\| = (\frac{1}{\dim(v)} \sum_{i=1}^{\dim(v)} |v_i|^p)^{1/p}$), then we simply use $N_p(\epsilon, \mathcal{F}, \mathcal{S})$ to denote the corresponding ℓ_p -covering number.

B. Proofs for Section 3.2

B.1. Proof of Lemma 3.5

Note that

$$\begin{aligned} & \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\ell(\hat{f}^X(X), Y) - \ell(f^*(X), Y) \right] \\ & \quad // \text{ We add and subtract loss of the local optimizer } f^{X,*}(\cdot) \text{ expected over } \mathcal{D}^{X,r} \\ & = \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\ell(\hat{f}^X(X), Y) - \mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} \left[\ell(f^{X,*}(X'), Y') \right] \right. \\ & \quad \left. + \mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} \left[\ell(f^{X,*}(X'), Y') \right] - \ell(f^*(X), Y) \right] \\ & \quad // \text{ We add and subtract loss of the global optimizer } f^*(\cdot) \text{ expected over } \mathcal{D}^{X,r} \\ & = \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\ell(\hat{f}^X(X), Y) - \mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} \left[\ell(f^{X,*}(X'), Y') \right] \right. \\ & \quad \left. + \mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} \left[\ell(f^*(X'), Y') \right] - \ell(f^*(X), Y) \right. \\ & \quad \left. + \mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} \left[\ell(f^{X,*}(X'), Y') \right] - \mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} \left[\ell(f^*(X'), Y') \right] \right] \\ & \quad // \text{ We group (1) local vs global optimizer, (2) global optimizer at } X \text{ vs expected over } \mathcal{D}^{X,r}, \\ & \quad // \text{ and (3) ERM loss at } X \text{ vs local optimizer loss expected over } \mathcal{D}^{X,r} \\ & = \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} \left[\ell(f^{X,*}(X'), Y') - \ell(f^*(X'), Y') \right] \right. \\ & \quad \left. + \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} \left[\ell(f^*(X'), Y') \right] - \ell(f^*(X), Y) \right] \right. \\ & \quad \left. + \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\ell(\hat{f}^X(X), Y) - \mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} \left[\ell(f^{X,*}(X'), Y') \right] \right] \right] \\ & \quad // \text{ We add and subtract loss of the empirical optimizer } \hat{f}^X(\cdot) \text{ expected over } \mathcal{D}^{X,r} \\ & = \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} \left[\ell(f^{X,*}(X'), Y') - \ell(f^*(X'), Y') \right] \right. \\ & \quad \left. + \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} \left[\ell(f^*(X'), Y') \right] - \ell(f^*(X), Y) \right] \right] \end{aligned}$$

$$\begin{aligned}
 & + \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\ell(\hat{f}^X(X), Y) - \mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} [\ell(\hat{f}^X(X'), Y')] \right. \\
 & \quad \left. + \mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} [\ell(\hat{f}^X(X'), Y')] - \mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} [\ell(f^{X,*}(X'), Y')] \right] \\
 // \text{ We (1) bound difference of loss at } X \text{ and loss expected over } \mathcal{D}^{X,r} \\
 & \quad \text{by maximizing over function class,} \\
 // \text{ and (2) subtract empirical loss of empirical optimizer and add (larger) empirical} \\
 & \quad \text{loss of local optimizer} \\
 \leq & \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} [\ell(f^{X,*}(X'), Y') - \ell(f^*(X'), Y')] \right] \\
 & + \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\sup_{f \in \mathcal{F}^{\text{global}}} \left| \mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} [\ell(f(X'), Y')] - \ell(f(X), Y) \right| \right] \\
 & + \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\sup_{f \in \mathcal{F}^{\text{loc}}} \left| \ell(f(X), Y) - \mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} [\ell(f(X'), Y')] \right| \right] \\
 & + \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} [\ell(\hat{f}^X(X'), Y')] - \frac{1}{|\mathcal{R}^X|} \sum_{(x',y') \in \mathcal{R}^X} \ell(\hat{f}^X(x'), y') \right] \\
 & + \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\frac{1}{|\mathcal{R}^X|} \sum_{(x',y') \in \mathcal{R}^X} \ell(f^{X,*}(x'), y') - \mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} [\ell(f^{X,*}(X'), Y')] \right] \tag{22}
 \end{aligned}$$

// We (1) bound difference of empirical vs expected loss of empirical optimizer
by maximizing over function class,

$$\begin{aligned}
 \leq & \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} [\ell(f^{X,*}(X'), Y') - \ell(f^*(X'), Y')] \right] \\
 & + \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\sup_{f \in \mathcal{F}^{\text{global}}} \left| \mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} [\ell(f(X'), Y')] - \ell(f(X), Y) \right| \right] \\
 & + \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\sup_{f \in \mathcal{F}^{\text{loc}}} \left| \ell(f(X), Y) - \mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} [\ell(f(X'), Y')] \right| \right] \\
 & + \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\sup_{f \in \mathcal{F}^{\text{loc}}} \left| \mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} [\ell(f(X'), Y')] - \frac{1}{|\mathcal{R}^X|} \sum_{(x',y') \in \mathcal{R}^X} \ell(f(x'), y') \right| \right] \\
 & + \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} [\ell(f^{X,*}(X'), Y')] - \frac{1}{|\mathcal{R}^X|} \sum_{(x',y') \in \mathcal{R}^X} \ell(f^{X,*}(x'), y') \right] \tag{23}
 \end{aligned}$$

□

B.2. Proof of Theorem 3.7

As discussed in Sec. 3, the proof of Theorem 3.7 requires bounding three terms in Lemma 3.5. We now proceed to establishing the desired bounds.

Local vs global loss. The local vs global loss can be bounded easily using the local regularity condition, and due to the fact that $\mathcal{F}^{\text{loc}} \approx \cup_x \mathcal{F}^x$. Let

$$f^{X,\text{loc}} = \arg \min_{f \in \mathcal{F}^X} \mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} [\ell(f(X'), Y')].$$

$$\begin{aligned}
 & \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} [\ell(f^{X,*}(X'), Y') - \ell(f^*(X'), Y')] \right] \\
 & \leq \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} [\ell(f^{X,*}(X'), Y') - \ell(f^{X,\text{loc}}(X'), Y')] \right] \\
 & \quad + \mathbb{E}_{(X,Y) \sim \mathcal{D}} \left[\mathbb{E}_{(X',Y') \sim \mathcal{D}^{X,r}} [\ell(f^{X,\text{loc}}(X'), Y') - \ell(f^*(X'), Y')] \right] \\
 & \leq \varepsilon_{\text{loc}} + \varepsilon_X.
 \end{aligned}$$

Global and local: Sample vs retrieved set risk. The following lemma bounds the second term in Lemma 3.5. Recall the definition, for any $L > 0$,

$$\mathcal{M}_r(L; \ell, f_{\text{true}}, \mathcal{F}) = 2L_\ell \left(Lr + (2\|\mathcal{F}\|_\infty - Lr) c_{\text{true}} (2L_{\text{true}}r)^{\alpha_{\text{true}}} \right). \quad (24)$$

Lemma B.1. *Under Assumption 3.1, for a L -coordinate Lipschitz function class \mathcal{F} with $\|\mathcal{F}\|_\infty := \sup_{x \in \mathcal{X}} \sup_{f \in \mathcal{F}} \|f(x)\|_\infty$ we have*

$$\begin{aligned} \mathbb{E}_{(X, Y) \sim \mathcal{D}} \left[\sup_{f \in \mathcal{F}} \left| \ell(f(X), Y) - \mathbb{E}_{(X', Y') \sim \mathcal{D}^{X, r}} [\ell(f(X'), Y')] \right| \right] \\ \leq 2L_\ell \left(Lr + (2\|\mathcal{F}\|_\infty - Lr) c_{\text{true}} (2L_{\text{true}}r)^{\alpha_{\text{true}}} \right). \end{aligned}$$

Proof. We are given the example (X, Y) . Let us fix an arbitrary $f \in \mathcal{F}$, and any arbitrary example (x', y') in the r neighborhood of X .

We first bound the perturbation in $\gamma_f(\cdot)$ for a given label \tilde{Y} .

$$\begin{aligned} |\gamma_f(X_1, \tilde{Y}) - \gamma_f(X_2, \tilde{Y})| &\leq |f_{\tilde{Y}}(X_1) - \max_{s \neq \tilde{Y}} f_s(X_1) - f_{\tilde{Y}}(X_2) + \max_{s' \neq \tilde{Y}} f_{s'}(X_2)| \\ &\leq |f_{\tilde{Y}}(X_1) - f_{\tilde{Y}}(X_2)| + \left| \max_{s \neq \tilde{Y}} f_s(X_1) - \max_{s' \neq \tilde{Y}} f_{s'}(X_2) \right| \\ &\leq |f_{\tilde{Y}}(X_1) - f_{\tilde{Y}}(X_2)| + \max_{s \neq \tilde{Y}} |f_s(X_1) - f_s(X_2)| \\ &\leq 2L \|X_1 - X_2\|_2 \end{aligned}$$

We can now proceed with bounding the loss.

$$\begin{aligned} |\ell(f(X), Y) - \ell(f(x'), y')| &= |\ell(\gamma_f(X, Y)) - \ell(\gamma_f(x', y'))| \\ &\leq L_\ell |\gamma_f(X, Y) - \gamma_f(x', y')| \\ &\leq \begin{cases} 4L_\ell \|f\|_\infty; Y \neq y' \\ 2L_\ell Lr; Y = y' \end{cases} \end{aligned}$$

Under Assumption 3.1, if we have $\gamma_{f_{\text{true}}}(X, Y) > 2L_{\text{true}}r$, then following the above argument we have $\gamma_{f_{\text{true}}}(X', Y) > 0$, thus Y is the true label of X' . In other words, $\gamma_{f_{\text{true}}}(X, Y) > 2L_{\text{true}}r$ imply for any X' in the r neighborhood of X its true label $Y' = Y$.

$$\begin{aligned} |\ell(f(X), Y) - \ell(f(x'), y')| \\ \leq 2L_\ell Lr \mathbb{1}(\gamma_{f_{\text{true}}}(X, Y) > 2L_{\text{true}}r) + 4L_\ell \|f\|_\infty \mathbb{1}(\gamma_{f_{\text{true}}}(X, Y) \leq 2L_{\text{true}}r) \\ \leq 2L_\ell Lr + 2L_\ell (2\|f\|_\infty - Lr) \mathbb{1}(\gamma_{f_{\text{true}}}(X, Y) \leq 2L_{\text{true}}r) \end{aligned}$$

As (x', y') was an arbitrary r -neighbor, we have

$$\begin{aligned} |\ell(f(X), Y) - \mathbb{E}_{(X', Y') \sim \mathcal{D}^{X, r}} \ell(f(X'), Y')| \\ \leq \mathbb{E}_{(X', Y') \sim \mathcal{D}^{X, r}} |\ell(f(X), Y) - \ell(f(X'), Y')| \\ \leq 2L_\ell Lr + 2L_\ell (2\|f\|_\infty - Lr) \mathbb{1}(\gamma_{f_{\text{true}}}(X, Y) \leq 2L_{\text{true}}r) \end{aligned}$$

Furthermore, as f was arbitrary, we have

$$\begin{aligned} \sup_{f \in \mathcal{F}} |\ell(f(X), Y) - \mathbb{E}_{(X', Y') \sim \mathcal{D}^{X, r}} \ell(f(X'), Y')| \\ \leq \sup_{f \in \mathcal{F}} 2L_\ell Lr + 2L_\ell (2\|f\|_\infty - Lr) \mathbb{1}(\gamma_{f_{\text{true}}}(X, Y) \leq 2L_{\text{true}}r) \end{aligned}$$

$$= 2L_\ell Lr + 2L_\ell(2\|\mathcal{F}\|_\infty - Lr) \mathbb{1}(\gamma_{f^{\text{true}}}(X, Y) \leq 2L_{\text{true}}r).$$

Note f^{true} is independent of f , which was used in the derivation of above inequalities. Taking expectation over (X, Y) , and using the margin condition as given in assumption 3.1 we obtain

$$\begin{aligned} & \mathbb{E}_{(X, Y) \sim \mathcal{D}} \left[\sup_{f \in \mathcal{F}} |\ell(f(X), Y) - \mathbb{E}_{(X', Y') \sim \mathcal{D}^{x, r}} \ell(f(X'), Y')| \right] \\ &= 2L_\ell Lr + 2L_\ell(2\|\mathcal{F}\|_\infty - Lr) \mathbb{P}_{(X, Y) \sim \mathcal{D}} \left[\gamma_{f^{\text{true}}}(X, Y) \leq 2L_{\text{true}}r \right] \\ &\leq 2L_\ell Lr + 2L_\ell(2\|\mathcal{F}\|_\infty - Lr) c_{\text{true}} (2L_{\text{true}}r)^{\alpha_{\text{true}}} = \mathcal{M}_r(L; \ell, f^{\text{true}}, \mathcal{F}). \end{aligned}$$

□

Plugging in the Lipschitz bounds for the function classes \mathcal{F}^{loc} and $\mathcal{F}^{\text{global}}$ in the above lemma bounds the second term.

An alternative way of bounding the risk difference is as follows:

$$\begin{aligned} & \mathbb{E}_{(X, Y) \sim \mathcal{D}} \left[\mathbb{E}_{(X', Y') \sim \mathcal{D}^{x, r}} [\ell(f(X'), Y')] - \ell(f(X), Y) \right] \\ &= \int_{x \in \mathcal{X}} \left(\int_{x' \in B(x, r) \cap \mathcal{X}} \frac{h(x') \rho_{\mathcal{D}}(x')}{\mathbb{P}_{\mathcal{D}}[B(x, r) \cap \mathcal{X}]} dx' \right) \rho_{\mathcal{D}}(x) dx - \int_{x \in \mathcal{X}} h(x) \rho_{\mathcal{D}}(x) dx \\ &= \int_{x' \in \mathcal{X}} h(x') \left(\int_{x \in B(x', r) \cap \mathcal{X}} \frac{\rho_{\mathcal{D}}(x)}{\mathbb{P}_{\mathcal{D}}[B(x, r) \cap \mathcal{X}]} dx \right) \rho_{\mathcal{D}}(x') dx' - \int_{x \in \mathcal{X}} h(x) \rho_{\mathcal{D}}(x) dx \\ &= \int_{x' \in \mathcal{X}} h(x') \left(\int_{x \in B(x', r) \cap \mathcal{X}} \frac{\rho_{\mathcal{D}}(x)}{\mathbb{P}_{\mathcal{D}}[B(x, r) \cap \mathcal{X}]} dx - 1 \right) \rho_{\mathcal{D}}(x') dx' \\ &\leq h_{\max} \int_{x' \in \mathcal{X}} \left| \int_{x \in B(x', r) \cap \mathcal{X}} \frac{\rho_{\mathcal{D}}(x)}{\mathbb{P}_{\mathcal{D}}[B(x, r) \cap \mathcal{X}]} dx - 1 \right| \rho_{\mathcal{D}}(x') dx' \\ &\leq h_{\max} \int_{x' \in \mathcal{X}} \max \left(\left| \frac{1}{1 \pm c_{\text{wdc}+} r^{\alpha_{\text{wdc}+}}} - 1 \right| \right) \rho_{\mathcal{D}}(x') dx' = \frac{h_{\max} c_{\text{wdc}+} r^{\alpha_{\text{wdc}+}}}{1 - c_{\text{wdc}+} r^{\alpha_{\text{wdc}+}}}. \end{aligned}$$

We can express $\ell(f(X), Y) = h(X)$ because Y is a deterministic function of X .

Under Assumption 3.4, with constants $c_{\text{wdc}+}$ and $\alpha_{\text{wdc}+}$, recall that

$$\left| \frac{\mathbb{P}_{X' \sim \mathcal{D}}[\text{dl}(X', x) \leq r]}{\rho_{\mathcal{D}}(x) \text{vol}_d(r)} - 1 \right| \leq c_{\text{wdc}+} r^{\alpha_{\text{wdc}+}}.$$

Plugging this in gives us the final inequality.

Example: Let us consider the term $\frac{\mathbb{P}_{\mathcal{D}}[B(x', r)]}{\mathbb{P}_{\mathcal{D}}[B(x, r) \cap \mathcal{X}]}$ for \mathcal{D} being multivariate Gaussian $N(\mu, \Sigma)$. Let $\text{vol}_d(r)$ imply the volume, and $S_d(r)$ the surface area of a d -sphere of radius r in dimension d .

$$\begin{aligned} \mathbb{P}_{\mathcal{D}}[B(x, r) \cap \mathcal{X}] &= \int_{z \in B(x, r) \cap \mathcal{X}} \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2}(z - \mu)^T \Sigma^{-1} (z - \mu)\right) dz \\ &= \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu)\right) \int_{z \in B(x, r) \cap \mathcal{X}} \exp\left(-\frac{1}{2}(z + x - 2\mu)^T \Sigma^{-1} (z - x)\right) dz \\ &= \rho_{\mathcal{D}}(x) \int_{u \in B(0, r)} \exp\left(-\frac{1}{2}(u + 2(x - \mu))^T \Sigma^{-1} u\right) du \\ &= \rho_{\mathcal{D}}(x) (\text{vol}_d(r) - c \int_{u \in B(0, r)} (u + 2(x - \mu))^T \Sigma^{-1} u du) \text{ for some } c \in [1/4, 1/2] \\ &= \rho_{\mathcal{D}}(x) (\text{vol}_d(r) - c \int_{u \in B(0, r)} (c_1 \|u\|_2^2 + 2(x - \mu)^T \Sigma^{-1} u) du) \text{ for some } c_1 \in [\lambda_{\min}(\Sigma^{-1}), \lambda_{\max}(\Sigma^{-1})] \end{aligned}$$

We have used $\exp(-x) \in (1 - x, 1 - x/2)$ for $x \leq 1.59$. Also, $\frac{u^T \Sigma^{-1} u}{\|u\|_2^2} \in [\lambda_{\min}(\Sigma^{-1}), \lambda_{\max}(\Sigma^{-1})]$.

Then using polar coordinate transform we obtain

$$\int_{u \in B(0,r)} \|u\|_2^2 du = \int_0^r l^2 S_d(l) dl = \int_0^r l^2 \frac{d l^{d-1} \pi^{d/2}}{\Gamma(1+d/2)} dl = r^{d+2} \frac{d \pi^{d/2}}{(d+2)\Gamma(1+d/2)} = \text{vol}_d(r) \frac{dr^2}{d+2}.$$

Let $\xi = 2(x - \mu)\Sigma^{-1}$. We want to integrate $\xi^T u$ over $B(0, r)$. Through a somewhat different polar transform where the polar axis is parallel to ξ and the angle of u and ξ is θ we can do the integral as follows for $d \geq 2$.

$$\int_{u \in B(0,r)} \xi^T u du = \int_{l=0}^r \int_{\theta=0}^{\pi} |\xi| l \cos(\theta) S_{d-1}(l) \sin^{d-2}(\theta) d\theta dl = \frac{d\pi^{d/2}}{\Gamma(1+d/2)} \int_{l=0}^r |\xi| l^{d-1} dl \underbrace{\int_{\theta=0}^{\pi} \cos(\theta) \sin^{d-2}(\theta) d\theta}_{=0} = 0.$$

Substituting, these values in the above inequality we get

$$\mathbb{P}_D[B(x, r) \cap \mathcal{X}] = \rho_D(x) \text{vol}_d(r) (1 - c_2 r^2), \text{ for some } c_2 = \left[\frac{d\lambda_{\min}(\Sigma^{-1})}{4(d+2)}, \frac{d\lambda_{\max}(\Sigma^{-1})}{2(d+2)} \right].$$

Therefore, the difference of Retrieved vs Sample risk for multi-variate Gaussian is bounded as

$$\frac{l_{\max}(f) dr^2}{(d+2)\lambda_{\max}(\Sigma)} \text{ for } r \leq \sqrt{(d+2)\lambda_{\max}(\Sigma)}.$$

Generalization of local ERM. Recall the function class $\mathcal{G}(X, Y) = \{\ell(\gamma_f(\cdot, \cdot)) - \ell(\gamma_f(X, Y)) : f \in \mathcal{F}^{\text{loc}}\}$. Here $\mathcal{G}(X, Y) : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$. Note that the function class is parameterized by (X, Y) . Let us define some quantities of the function class on a set $S \subseteq \mathcal{X} \times \mathcal{Y}$ as

$$\mathcal{G}_{\max}((X, Y); S) = \sup_{g \in \mathcal{G}(X, Y)} \sup_{(x', y') \in S} |g(x', y')|$$

By centering each function $f \in \mathcal{F}^{\text{loc}}$ at the point (X, Y) we can transform the generalization over the function class \mathcal{F}^{loc} , to the generalization over the function class $\mathcal{G}(X, Y)$. In particular, we have

$$\begin{aligned} & \mathbb{E}_{(X, Y) \sim D} \left[\sup_{f \in \mathcal{F}^{\text{loc}}} \left| \mathbb{E}_{(X', Y') \sim D^{X, r}} [\ell(f(X'), Y')] - \frac{1}{|\mathcal{R}^X|} \sum_{(x', y') \in \mathcal{R}^X} \ell(f(x'), y') \right| \right] \\ & \leq \mathbb{E}_{(X, Y) \sim D} \left[\sup_{f \in \mathcal{F}^{\text{loc}}} \left| \mathbb{E}_{(X', Y') \sim D^{X, r}} [\ell(f(X'), Y') - \ell(f(X), Y)] \right. \right. \\ & \quad \left. \left. - \frac{1}{|\mathcal{R}^X|} \sum_{(x', y') \in \mathcal{R}^X} \ell(f(x'), y') - \ell(f(X), Y) \right| \Big| |\mathcal{R}^X| \geq N(r, \delta) \right] + 4\delta L_\ell \|\mathcal{F}^{\text{loc}}\|_\infty \\ & = \mathbb{E}_{(X, Y) \sim D} \left[\sup_{g \in \mathcal{G}(X, Y)} \left| \mathbb{E}_{(X', Y') \sim D^{X, r}} [g(X', Y')] - \frac{1}{|\mathcal{R}^X|} \sum_{(x', y') \in \mathcal{R}^X} g(x', y') \right| \Big| |\mathcal{R}^X| \geq N(r, \delta) \right] + 4\delta L_\ell \|\mathcal{F}^{\text{loc}}\|_\infty. \end{aligned}$$

We next state a standard result of learning theory that bounds the final term using the Rademacher complexity of the function class $\mathcal{G}(X, Y)$ (Shalev-Shwartz & Ben-David, 2014).

Lemma B.2 (Adapted from Theorem 26.5 in Shalev-Shwartz & Ben-David (2014)). *For any $(X, Y) \in \mathcal{X} \times \mathcal{Y}$ and a neighborhood set \mathcal{R}^X , and any function $g \in \mathcal{G}(X, Y)$, for each $\delta > 0$ with probability at least $(1 - \delta)$ the following holds*

$$\left| \mathbb{E}_{(X', Y') \sim D^{X, r}} [g(X', Y')] - \frac{1}{|\mathcal{R}^X|} \sum_{(x', y') \in \mathcal{R}^X} g(x', y') \right| \leq 2\mathfrak{R}_{\mathcal{R}^X}(\mathcal{G}(X, Y)) + 4\mathcal{G}_{\max}((X, Y); \mathcal{R}^X) \sqrt{\frac{2 \ln(4/\delta)}{|\mathcal{R}^X|}}.$$

Taking expectation with respect to (X, Y) , we obtain

$$\mathbb{E}_{(X, Y) \sim D} \left[\sup_{g \in \mathcal{G}(X, Y)} \left| \mathbb{E}_{(X', Y') \sim D^{X, r}} [g(X', Y')] - \frac{1}{|\mathcal{R}^X|} \sum_{(x', y') \in \mathcal{R}^X} g(x', y') \right| \Big| |\mathcal{R}^X| \geq N(r, \delta) \right]$$

$$\begin{aligned}
 &\leq 2\mathbb{E}_{(X,Y)\sim\mathcal{D}}\left[\mathfrak{R}_{\mathcal{R}^X}(\mathcal{G}(X,Y))\Big|\mathcal{R}^X\geq N(r,\delta)\right]+ \\
 &\quad 4\mathbb{E}_{(X,Y)\sim\mathcal{D}}\left[\mathcal{G}_{\max}((X,Y);\mathcal{R}^X)\sqrt{\frac{2\ln(4/\delta)}{|\mathcal{R}^X|}}\Big|\mathcal{R}^X\geq N(r,\delta)\right]+4\delta L_\ell\|\mathcal{F}^{\text{loc}}\|_\infty \\
 &\leq 2\mathbb{E}_{(X,Y)\sim\mathcal{D}}\left[\mathfrak{R}_{\mathcal{R}^X}(\mathcal{G}(X,Y))\Big|\mathcal{R}^X\geq N(r,\delta)\right]+ \\
 &\quad 4\mathbb{E}_{(X,Y)\sim\mathcal{D}}\left[\mathcal{G}_{\max}((X,Y);\mathcal{R}^X)\Big|\mathcal{R}^X\geq N(r,\delta)\right]\mathbb{E}_{(X,Y)\sim\mathcal{D}}\left[\sqrt{\frac{2\ln(4/\delta)}{|\mathcal{R}^X|}}\Big|\mathcal{R}^X\geq N(r,\delta)\right]+4\delta L_\ell\|\mathcal{F}^{\text{loc}}\|_\infty \\
 &\leq 2\mathbb{E}_{(X,Y)\sim\mathcal{D}}\left[\mathfrak{R}_{\mathcal{R}^X}(\mathcal{G}(X,Y))\Big|\mathcal{R}^X\geq N(r,\delta)\right]+4\mathcal{M}_r(L_{\text{loc}};\ell,f_{\text{true}},\mathcal{F}^{\text{loc}})\mathbb{E}_{(X,Y)\sim\mathcal{D}}\left[\sqrt{\frac{2\ln(4/\delta)}{N(r,\delta)}}\right]+4\delta L_\ell\|\mathcal{F}^{\text{loc}}\|_\infty.
 \end{aligned}$$

In the first inequality, with probability $(1-\delta)$ we apply the bound from Lemma B.2, whereas we use the bound $4L_\ell\|\mathcal{F}^{\text{loc}}\|_\infty$ with remaining probability δ . Also from the proof of Lemma B.1 we have that

$$\mathcal{G}_{\max}((X,Y);\mathcal{R}^X)\leq 2L_\ell\left(Lr+(\max\{Lr,2\|\mathcal{F}^{\text{loc}}\|_\infty\}-Lr)\mathbb{1}(\gamma_{f_{\text{true}}}(X,Y)\leq 2L_{\text{true}}r)\right).$$

Taking expectation with respect to \mathcal{D} completes the bound. While taking expectation we crucially use the fact that $\gamma_{f_{\text{true}}}(X,Y)$ is independent of $|\mathcal{R}^X|$ to arrive at the $\mathcal{M}_r(L_{\text{loc}};\ell,f_{\text{true}},\mathcal{F}^{\text{loc}})$ bound.

Central absolute moment of $f^{X,*}$. As the function $f^{X,*}$ is fixed using centering, and then Hoeffding bound, we can directly bound the remaining term. We have with probability at least $(1-\delta)$

$$\begin{aligned}
 &\left|\mathbb{E}_{(X',Y')\sim\mathcal{D}^{X,r}}[\ell(f^{X,*}(X'),Y')]-\frac{1}{|\mathcal{R}^X|}\sum_{(x',y')\in\mathcal{R}^X}\ell(f^{X,*}(x'),y')\right| \\
 &= \left|\mathbb{E}_{(X',Y')\sim\mathcal{D}^{X,r}}[\ell(f^{X,*}(X'),Y')-\ell(f^{X,*}(X),Y)]-\frac{1}{|\mathcal{R}^X|}\sum_{(x',y')\in\mathcal{R}^X}\ell(f^{X,*}(x'),y')-\ell(f^{X,*}(X),Y)\right| \\
 &\leq \mathcal{G}_{\max}((X,Y);\mathcal{R}^X)\sqrt{\frac{\ln(2/\delta)}{|\mathcal{R}^X|}}
 \end{aligned}$$

Taking expectation similar to the previous case we obtain,

$$\begin{aligned}
 &\mathbb{E}_{(X,Y)\sim\mathcal{D}}\left[\left|\mathbb{E}_{(X',Y')\sim\mathcal{D}^{X,r}}[\ell(f^{X,*}(X'),Y')]-\frac{1}{|\mathcal{R}^X|}\sum_{(x',y')\in\mathcal{R}^X}\ell(f^{X,*}(x'),y')\right|\right] \\
 &\leq \mathbb{E}_{(X,Y)\sim\mathcal{D}}\left[\min\{4L_\ell\|\mathcal{F}^{\text{loc}}\|_\infty,\mathcal{G}_{\max}((X,Y);\mathcal{R}^X)\sqrt{\frac{\ln(2/\delta)}{|\mathcal{R}^X|}}\}\right] \\
 &\leq \mathcal{M}_r(L_{\text{loc}};\ell,f_{\text{true}},\mathcal{F}^{\text{loc}})\mathbb{E}_{(X,Y)\sim\mathcal{D}}\left[\sqrt{\frac{\ln(2/\delta)}{N(r,\delta)}}\right]+4\delta L_\ell\|\mathcal{F}^{\text{loc}}\|_\infty.
 \end{aligned}$$

Here, we use the fact that $\gamma_{f_{\text{true}}}(X,Y)$ is independent of $|\mathcal{R}^X|$. This concludes the proof of Theorem 3.7.

B.3. Bounding the Rademacher complexity $\mathfrak{R}_{\mathcal{R}^X}(\mathcal{G}(X,Y))$

We now derive bounds on the Rademacher complexity of the class $\mathcal{G}(X,Y)$. We use the covering number based bounds for that purpose. We then start by relating it to the covering number of the \mathcal{F}^{loc} function class. Finally, we provide a bound on the class of functions residing in bounded norm Reproducing Kernel Hilbert Space.

We will use $\mathcal{G}_{\max}(X,Y)$ instead of $\mathcal{G}_{\max}((X,Y);\mathcal{R}^X)$ when the context is clear. Similar to $\mathcal{G}(X,Y)$, we define the function class $\mathcal{G}=\{\ell(\gamma_f(\cdot,\cdot)):f\in\mathcal{F}^{\text{loc}}\}$ which does not depend on the locality centered around (X,Y) . On a set $S\subseteq\mathcal{X}\times\mathcal{Y}$ we can define $\mathcal{G}_{\max}(S)=\sup_{g\in\mathcal{G}}\sup_{(x',y')\in S}|g(x',y')|$.

Lemma B.3. *Under Assumption 3.1 we have for any retrieved set within radius r of X , \mathcal{R}^X , for any $p\geq 1$*

$$\mathfrak{R}_{\mathcal{R}^X}(\mathcal{G}(X,Y))$$

$$\leq \begin{cases} \mathcal{G}_{\max}(X, Y) \\ \inf_{\epsilon \in [0, \mathcal{G}_{p, \max}(X, Y)/2]} \left(4\epsilon + \frac{12}{\sqrt{|\mathcal{R}^X|}} \int_{\epsilon}^{\mathcal{G}_{p, \max}(X, Y)/2} \sqrt{\log\left(\frac{2\mathcal{G}_{\max}}{\nu}\right) \log\left(\mathcal{N}_p(\nu/2, \mathcal{G}, \mathcal{R}^X)\right)} d\nu \right) \\ \inf_{\epsilon \in [0, \mathcal{G}_{\max}(X, Y)/2]} \left(4\epsilon + \frac{12}{\sqrt{|\mathcal{R}^X|}} \int_{\epsilon}^{\mathcal{G}_{\max}(X, Y)/2} \sqrt{\log\left(\mathcal{N}_{\infty}(\nu/2, \mathcal{G}, \mathcal{R}^X \cup \{(X, Y)\})\right)} d\nu \right). \end{cases}$$

As a corollary we obtain the following rates, as the log-covering number varies with ν at different rates.

Corollary B.4. *Under Assumption 3.1 we have for any retrieved set within radius r of X , \mathcal{R}^X , for any $p \geq 1$*

$$\mathfrak{R}_{\mathcal{R}^X}(\mathcal{G}(X, Y)) \leq \begin{cases} C'(p, \mathcal{F}^{\text{loc}}) \frac{\log^2(2 \max\{|\mathcal{R}^X|, \mathcal{G}_{\max}\})}{\sqrt{|\mathcal{R}^X|}}; & \text{if } \log\left(\mathcal{N}_p(\epsilon, \mathcal{G}, n)\right) \leq C^2(p, \mathcal{F}^{\text{loc}}) \log(n/\epsilon)/\epsilon^2, \\ \frac{C'(p, \mathcal{F}^{\text{loc}}) \mathcal{G}_{p, \max}(X, Y)^{1-\alpha/2} \log(2 \max\{|\mathcal{R}^X|, \mathcal{G}_{\max}\})}{\sqrt{|\mathcal{R}^X|}}; & \text{if } \log\left(\mathcal{N}_p(\epsilon, \mathcal{G}, n)\right) \leq C^2(p, \mathcal{F}^{\text{loc}}) \log(n/\epsilon)/\epsilon^{\alpha}, \alpha \in [0, 2). \end{cases}$$

Proof. Case $\log\left(\mathcal{N}_p(\epsilon, \mathcal{G}, n)\right) \leq C^2(p, \mathcal{F}^{\text{loc}}) \log(n/\epsilon)/\epsilon^{\alpha}$, $\alpha \in [0, 2)$:

$$\begin{aligned} \mathfrak{R}_{\mathcal{R}^X}(\mathcal{G}(X, Y)) &\leq \frac{4\mathcal{G}_{p, \max}(X, Y)}{\sqrt{|\mathcal{R}^X|}} + C' \frac{C(p, \mathcal{F}^{\text{loc}})}{\sqrt{|\mathcal{R}^X|}} \int_{\mathcal{G}_{p, \max}(X, Y)/\sqrt{|\mathcal{R}^X|}}^{\mathcal{G}_{p, \max}(X, Y)/2} \sqrt{\log\left(\frac{2\mathcal{G}_{\max}}{\nu}\right) \log\left(\frac{2|\mathcal{R}^X|}{\nu}\right)} \nu^{-\alpha/2} d\nu \\ &\leq \frac{4\mathcal{G}_{p, \max}(X, Y)}{\sqrt{|\mathcal{R}^X|}} + C' \frac{C(p, \mathcal{F}^{\text{loc}})}{\sqrt{|\mathcal{R}^X|}} \int_{\mathcal{G}_{p, \max}(X, Y)/\sqrt{|\mathcal{R}^X|}}^{\mathcal{G}_{p, \max}(X, Y)/2} \log\left(\frac{2 \max\{|\mathcal{R}^X|, \mathcal{G}_{\max}\}}{\nu}\right) \nu^{-\alpha/2} d\nu \\ &\leq \frac{4\mathcal{G}_{p, \max}(X, Y)}{\sqrt{|\mathcal{R}^X|}} + C'' \frac{C(p, \mathcal{F}^{\text{loc}})}{\sqrt{|\mathcal{R}^X|}} \left(\mathcal{G}_{p, \max}(X, Y)^{1-\alpha/2} \log(2 \max\{|\mathcal{R}^X|, \mathcal{G}_{\max}\}) + (1 - \alpha/2)^{-1} e^{-1} \right). \end{aligned}$$

The last inequality follows from

$$\begin{aligned} (1 - \alpha/2)^2 \int_c^b x^{-\alpha/2} \log(a/x) dx &= b^{1-\alpha/2} (1 + (1 - \alpha/2) \log(a/b)) - c^{1-\alpha/2} (1 + (1 - \alpha/2) \log(a/b)) \\ &= (b^{1-\alpha/2} - c^{1-\alpha/2}) (1 + (1 - \alpha/2) \log(a)) + b^{1-\alpha/2} \log(1/b) - c^{1-\alpha/2} \log(1/c) \\ &\leq (b^{1-\alpha/2} - c^{1-\alpha/2}) (1 + (1 - \alpha/2) \log(a)) + (1 - \alpha/2)^{-1} e^{-1} \end{aligned}$$

Case $\log\left(\mathcal{N}_p(\epsilon, \mathcal{G}, n)\right) \leq C^2(p, \mathcal{F}^{\text{loc}}) \log(n/\epsilon)/\epsilon^2$:

$$\begin{aligned} \mathfrak{R}_{\mathcal{R}^X}(\mathcal{G}(X, Y)) &\leq \frac{4}{\sqrt{|\mathcal{R}^X|}} + C' \frac{C(p, \mathcal{F}^{\text{loc}})}{\sqrt{|\mathcal{R}^X|}} \int_{1/\sqrt{|\mathcal{R}^X|}}^{\mathcal{G}_{p, \max}(X, Y)/2} \sqrt{\log\left(\frac{2\mathcal{G}_{\max}}{\nu}\right) \log\left(\frac{2|\mathcal{R}^X|}{\nu}\right)} \nu^{-1} d\nu \\ &\leq \frac{4}{\sqrt{|\mathcal{R}^X|}} + C' \frac{C(p, \mathcal{F}^{\text{loc}})}{\sqrt{|\mathcal{R}^X|}} \int_{1/\sqrt{|\mathcal{R}^X|}}^{\mathcal{G}_{p, \max}(X, Y)/2} \log\left(\frac{2 \max\{|\mathcal{R}^X|, \mathcal{G}_{\max}\}}{\nu}\right) \nu^{-1} d\nu \\ &\leq \frac{4}{\sqrt{|\mathcal{R}^X|}} + C' \frac{C(p, \mathcal{F}^{\text{loc}})}{\sqrt{|\mathcal{R}^X|}} \left(\log^2(2 \max\{|\mathcal{R}^X|, \mathcal{G}_{\max}\}) \sqrt{|\mathcal{R}^X|} - \log^2\left(\frac{4 \max\{|\mathcal{R}^X|, \mathcal{G}_{\max}\}}{\mathcal{G}_{p, \max}(X, Y)}\right) \right) \\ &\leq \frac{4}{\sqrt{|\mathcal{R}^X|}} + C' \frac{C(p, \mathcal{F}^{\text{loc}})}{\sqrt{|\mathcal{R}^X|}} \log^2(2 \max\{|\mathcal{R}^X|, \mathcal{G}_{\max}\}) \sqrt{|\mathcal{R}^X|}. \end{aligned}$$

□

Proof of Lemma B.3. Given the set \mathcal{R}^X , and some function $g \in \mathcal{G}(X, Y)$ let us define for $p \geq 1$

$$\|g\|_{p, \mathcal{R}^X} = \left(\frac{1}{|\mathcal{R}^X|} \sum_{(x', y') \in \mathcal{R}^X} |g(x', y')|^p \right)^{1/p}.$$

Then, we have $\mathcal{G}_{p,\max}((X, Y); \mathcal{R}^X) = \max_{g \in \mathcal{G}} \|g\|_{p, \mathcal{R}^X}$ for all $g \in \mathcal{G}(X, Y)$. For the sake of brevity we will use $\mathcal{G}_{p,\max}(X, Y)$ in place of $\mathcal{G}_{p,\max}((X, Y); \mathcal{R}^X)$. Note that we have from previous definition $\mathcal{G}_{\max}(X, Y) = \mathcal{G}_{\infty,\max}(X, Y) \geq \mathcal{G}_{p,\max}(X, Y)$ for any $p \geq 1$.

A simple bound on the Rademacher complexity comes as a function of the radius r but which is independent of the size of $|\mathcal{R}^X|$. Specifically, we have for Rademacher random variable σ_i -s

$$\mathfrak{R}_{\mathcal{R}^X}(\mathcal{G}(X, Y)) \leq \frac{1}{|\mathcal{R}^X|} \mathbb{E}_\sigma \left[\left| \sum_{(X'_i, Y'_i) \in \mathcal{R}^X} \sigma_i g_i \right| \right] \leq \max_{g_i \in \mathcal{G}(X, Y)} |g_i| \leq \mathcal{G}_{\max}(X, Y).$$

Next using the Chaining method (Shalev-Shwartz & Ben-David, 2014, Chapter 27) we can bound the Rademacher complexity as

$$\mathfrak{R}_{\mathcal{R}^X}(\mathcal{G}(X, Y)) \leq \inf_{\epsilon \in [0, \mathcal{G}_{p,\max}(X, Y)/2]} \left(4\epsilon + \frac{12}{\sqrt{|\mathcal{R}^X|}} \int_\epsilon^{\mathcal{G}_{p,\max}(X, Y)/2} \sqrt{\log \mathcal{N}_p(\nu, \mathcal{G}(X, Y), \mathcal{R}^X)} d\nu \right).$$

To finish the proof we need to show, for $p \geq 1$

$$\mathcal{N}_p(\nu, \mathcal{G}(X, Y), \mathcal{R}^X) \leq \mathcal{N}_p(\nu/2, \mathcal{G}, \mathcal{R}^X) \mathcal{N}_p(\nu/2, \mathcal{G}, \{(X, Y)\}).$$

First we fix any $p \geq 1$. Let $\widehat{\mathcal{U}}$ (a set of real numbers) be a $\nu/2$ cover (in ℓ_p norm) of \mathcal{G} with respect to $\{(X, Y)\}$. We have $\mathcal{N}_p(\nu, \mathcal{G}(X, Y), \mathcal{R}^X) \leq \frac{2\mathcal{G}_{\max}}{\nu}$ for any $p \geq 1$ and any $\nu > 0$. Further, let $\tilde{\mathcal{U}}$ be a $\nu/2$ cover of \mathcal{G} with respect to \mathcal{R}^X . Note for any $\tilde{u} \in \tilde{\mathcal{U}}$ we have $\tilde{u} \in \mathbb{R}^{|\mathcal{R}^X|}$.

Now, we fix any $g' \in \mathcal{G}$. We have at least one $\tilde{u} \in \tilde{\mathcal{U}}$, and $\hat{u} \in \widehat{\mathcal{U}}$ such that

$$\left(\frac{1}{|\mathcal{R}^X|} \sum_{(x', y') \in \mathcal{R}^X} |g'(x', y') - \tilde{u}(x', y')|^p \right)^{1/p} \leq \nu/2, \text{ and } |g'(X, Y) - \hat{u}| \leq \nu/2.$$

Therefore,

$$\begin{aligned} & \left(\frac{1}{|\mathcal{R}^X|} \sum_{(x', y') \in \mathcal{R}^X} |(g'(x', y') - g'(X, Y)) - (\tilde{u}(x', y') - \hat{u})|^p \right)^{1/p} \\ &= \left(\frac{1}{|\mathcal{R}^X|} \sum_{(x', y') \in \mathcal{R}^X} |(g'(x', y') - \tilde{u}(x', y')) + (\hat{u} - g'(X, Y))|^p \right)^{1/p} \\ &\leq \left(\frac{1}{|\mathcal{R}^X|} \sum_{(x', y') \in \mathcal{R}^X} |g'(x', y') - \tilde{u}(x', y')|^p \right)^{1/p} + |\hat{u} - g'(X, Y)| \\ &\leq \nu/2 + \nu/2 \leq \nu \end{aligned}$$

The first inequality follows by applying Minkowski's inequality. Whereas, for the second inequality we apply Jensen's inequality for $(\cdot)^{1/p}$ being a concave function for $p \geq 1$, and applying the appropriate scaling. Therefore, given the covers $\tilde{\mathcal{U}}$ and $\widehat{\mathcal{U}}$, we can construct the set \mathcal{U}' with entries $u' \in \mathbb{R}^{|\mathcal{R}^X|}$ as: $\mathcal{U}' := \{u' = (\tilde{u}(x, y) - \hat{u}) : \tilde{u} \in \tilde{\mathcal{U}}, \hat{u} \in \widehat{\mathcal{U}}\}$. In particular, $|\mathcal{U}'| = |\tilde{\mathcal{U}}| |\widehat{\mathcal{U}}|$. As the choice of $g' \in \mathcal{G}$ and $(x', y') \in \mathcal{R}^X$ were arbitrary, we have \mathcal{U}' to be the cover of $\mathcal{G}(X, Y)$.

For $p = \infty$ we can specialize the bound. In particular, consider \mathcal{U} to be a $\nu/2$ cover (in ℓ_∞ norm) of \mathcal{G} with respect to $\mathcal{R}^X \cup \{(X, Y)\}$. Then $\mathcal{U}' := \{u' = (\tilde{u}(x, y) - \hat{u}(X, Y)) : \tilde{u} \in \mathcal{U}\}$ creates a (normalized) ℓ_∞ cover for \mathcal{G} with respect to \mathcal{R}^X . This is true because $\left(\frac{1}{|\mathcal{R}^X|} \sum_{(x', y') \in \mathcal{R}^X} |g'(x', y') - \tilde{u}(x', y')|^p \right)^{1/p} \leq |g' - \tilde{u}|_\infty = \nu/2$ and $|\hat{u} - g'(X, Y)| \leq |g' - \tilde{u}|_\infty = \nu/2$. This concludes the proof. \square

The first term in the above Lemma is similar to the Chaining based Rademacher bounds (Shalev-Shwartz & Ben-David, 2014, Chapter 28) for \mathcal{G} , but the ϵ (in inf and in the integral) varies in $[0, \mathcal{G}_{\max}(X, Y)]$ instead of $[0, \mathcal{G}_{\max}]$. For small r we have $\mathcal{G}_{\max}(X, Y) \ll \mathcal{G}_{\max}$, which can be leveraged to give tight bounds in certain situations.

Example: $\mathcal{F}^{\text{loc}} \equiv \ell_\infty$ -bounded RKHS (Zhang, 2004): Let us consider the setting of Zhang (2004). In this setting, given some Reproducing Kernel Hilbert Space (RKHS) H , and a function $\tilde{f} \in H$, we can define the function $\tilde{f}(\cdot) = \tilde{f} \circ h_x$ where for some $h \in H$. We further define the set of functions with bounded norm

$$H_A = \{\tilde{f}(\cdot) \in H : \|\tilde{f}\|_H \sup_{x \in \mathcal{X}} \|h_x\|_H \leq A\}.$$

Finally, our local function class can be defined as

$$\mathcal{F}^{\text{loc}} = H_A^{|\mathcal{Y}|} = \{f(\cdot) : f_y(\cdot) \in H_A, \forall y \in \mathcal{Y}\}.$$

We have $\|\mathcal{F}^{\text{loc}}\|_\infty = A$. Recall that loss function for any $y \in \mathcal{Y}$ is given as $\ell(\gamma_f(x, y))$, for any $f \in \mathcal{F}^{\text{loc}}$. We also have for all $y \in \mathcal{Y}$, $|\ell(\gamma_f(x, y)) - \ell(\gamma_{f'}(x, y))| \leq 2L_\ell \sup_y |f_y(x) - f'_y(x)|$ (Zhang, 2004, Assumption 15) with $\gamma_A = 2L_\ell$.

Given the above setting, following Lemma 17 in Zhang (2004)³, we have for a universal constant c

$$\log \left(\mathcal{N}_\infty(2L_\ell\nu, \mathcal{G}, \mathcal{R}^X \cup \{(X, Y)\}) \right) \leq c|\mathcal{Y}| \|\mathcal{F}^{\text{loc}}\|_\infty^2 \frac{\ln(2 + \|\mathcal{F}^{\text{loc}}\|_\infty/\nu) + \ln(|\mathcal{R}^X| + 1)}{\nu^2}.$$

This gives us the following bound for the Rademacher complexity of \mathcal{F}^{loc}

$$\mathfrak{R}_{\mathcal{R}^X} \leq O\left(\sqrt{|\mathcal{Y}|} L_\ell \|\mathcal{F}^{\text{loc}}\|_\infty \frac{\ln(|\mathcal{R}^X| + 1)^{3/2}}{\sqrt{|\mathcal{R}^X|}}\right). \quad (25)$$

Proof of Equation (25). Without optimizing over ϵ above, we plug in $\epsilon = \frac{\mathcal{G}_{\max}(X, Y)}{\sqrt{|\mathcal{R}^X|}}$. We obtain

$$\begin{aligned} & \mathfrak{R}_{\mathcal{R}^X}(\mathcal{G}(X, Y)) \\ & \leq \frac{4\mathcal{G}_{\max}(X, Y)}{\sqrt{|\mathcal{R}^X|}} + \frac{12}{\sqrt{|\mathcal{R}^X|}} \int_{\frac{\mathcal{G}_{\max}(X, Y)}{\sqrt{|\mathcal{R}^X|}}}^{\mathcal{G}_{\max}(X, Y)/2} \sqrt{\log \left(\mathcal{N}_\infty \left(\nu/2, \mathcal{G}, \mathcal{R}^X \cup \{(X, Y)\} \right) \right)} d\nu \\ & \leq \frac{4\mathcal{G}_{\max}(X, Y)}{\sqrt{|\mathcal{R}^X|}} + \frac{48\sqrt{c|\mathcal{Y}|} L_\ell \|\mathcal{F}^{\text{loc}}\|_\infty}{\sqrt{|\mathcal{R}^X|}} \int_{\frac{\mathcal{G}_{\max}(X, Y)}{\sqrt{|\mathcal{R}^X|}}}^{\mathcal{G}_{\max}(X, Y)/2} \sqrt{\frac{\ln(2 + 4L_\ell \|\mathcal{F}^{\text{loc}}\|_\infty/\nu) + \ln(|\mathcal{R}^X| + 1)}{\nu^2}} d\nu \\ & \leq \frac{4\mathcal{G}_{\max}(X, Y)}{\sqrt{|\mathcal{R}^X|}} + \frac{48\sqrt{c|\mathcal{Y}|} L_\ell \|\mathcal{F}^{\text{loc}}\|_\infty}{\sqrt{|\mathcal{R}^X|}} \int_{\frac{\mathcal{G}_{\max}(X, Y)}{\sqrt{|\mathcal{R}^X|}}}^{\mathcal{G}_{\max}(X, Y)/2} \sqrt{\frac{\ln((\mathcal{G}_{\max}(X, Y) + 4L_\ell \|\mathcal{F}^{\text{loc}}\|_\infty)/\nu) + \ln(|\mathcal{R}^X| + 1)}{\nu^2}} d\nu \\ & \leq \frac{4\mathcal{G}_{\max}(X, Y)}{\sqrt{|\mathcal{R}^X|}} + \frac{48\sqrt{c|\mathcal{Y}|} L_\ell \|\mathcal{F}^{\text{loc}}\|_\infty}{\sqrt{|\mathcal{R}^X|}} \int_{\frac{1}{\sqrt{|\mathcal{R}^X|}}}^{1/2} \sqrt{\frac{\ln((1 + 4L_\ell \|\mathcal{F}^{\text{loc}}\|_\infty/\mathcal{G}_{\max}(X, Y))/\nu') + \ln(|\mathcal{R}^X| + 1)}{\nu'^2}} d\nu' \\ & \leq \frac{4\mathcal{G}_{\max}(X, Y)}{\sqrt{|\mathcal{R}^X|}} + \frac{32\sqrt{c|\mathcal{Y}|} L_\ell \|\mathcal{F}^{\text{loc}}\|_\infty}{\sqrt{|\mathcal{R}^X|}} \left(\ln \left((1 + 4L_\ell \|\mathcal{F}^{\text{loc}}\|_\infty/\mathcal{G}_{\max}(X, Y)) \sqrt{|\mathcal{R}^X|} \right) + \ln(|\mathcal{R}^X| + 1) \right)^{3/2} \end{aligned}$$

We use $\int_x \sqrt{\ln(a/x) + b/x} dx = -2/3(\ln(a/x) + b)^{3/2}$ for the final inequality, and ignore the negative part. \square

Example: $\mathcal{F}^{\text{loc}} \equiv \ell_2$ bounded RKHS (Lei et al., 2019): We consider a fixed kernel $K(x, x') = \langle \phi(x), \phi(x') \rangle$ for $x, x' \in \mathcal{X}$, and let H_K be the RKHS induced by K . Let us define the $\ell_{p, q}$ norm for the vectors $W = (w_1, w_2, \dots, w_{|\mathcal{Y}|}) \in H_K^{|\mathcal{Y}|}$ as $\|(w_1, \dots, w_{|\mathcal{Y}|})\|_{p, q} = \|(\|w_1\|_p, \dots, \|w_{|\mathcal{Y}|}\|_p)\|_q$.

For some norm bound $\Lambda > 0$, the local hypothesis space is defined as

$$\mathcal{F}^{\text{loc}} = \{f(\cdot) : f_y(\cdot) = \langle w_y, \phi(\cdot) \rangle, w_y \in H_K, \forall y \in \mathcal{Y}, \|(w_1, \dots, w_{|\mathcal{Y}|})\|_{2, 2} \leq \Lambda\}.$$

Recall that we have the loss function class $\mathcal{G} = \{\ell(\gamma_f(\cdot, \cdot)) : f \in \mathcal{F}^{\text{loc}}\}$, where the loss function $\ell(\cdot)$ is assumed to be L -Lipschitz continuous w.r.t. ℓ_∞ norm.

³We correct for a typographical error in Zhang (2004), where the $n \equiv |\mathcal{R}^X|$ comes in the denominator of the bound presented in Lemma 17. But Theorem 4 of Zhang (2002) shows this is a typographical error. Indeed, the covering number is not supposed to decrease with increasing number of points.

Given the retrieved set \mathcal{R}^X for some positive integer $n \geq 1$, $\tilde{\mathcal{F}}^X$ after Equation (8) in Lei et al. (2019) induced by \mathcal{R}^X .⁴ Let the worst case Rademacher complexity of a function class \mathcal{F} over n points be defined as $\mathfrak{R}_n(\mathcal{F})$. Also, for a set S let $\hat{B}(S) = \max_{(x,y) \in S} \sup_{W: \|W\|_{2,2} \leq \Lambda} \langle w_y, \phi(x) \rangle$. We have from Theorem 23 in Lei et al. (2019) that the covering number is bounded as follows: for any set $S = \{(x_i, y_i) : i = 1, \dots, n\}$ of size $n \geq 1$, for any $\varepsilon > 4L\mathfrak{R}_{n|y|}(\tilde{\mathcal{F}}^X)$

$$\log \left(\mathcal{N}_\infty(\varepsilon, \mathcal{G}, S) \right) \leq \frac{16n|y|L^2(\mathfrak{R}_{n|y|}(\tilde{\mathcal{F}}^X))^2}{\varepsilon^2} \log \left(\frac{2en|y|\hat{B}(S)L}{\varepsilon} \right).$$

Furthermore, from equation (18) in Lei et al. (2019) we have for any set

$$\frac{\Lambda \max_{(x,y) \in S} \|\phi(x)\|_2}{\sqrt{2n|y|}} \leq \mathfrak{R}_{n|y|}(\tilde{\mathcal{F}}^X) \leq \frac{\Lambda \max_{(x,y) \in S} \|\phi(x)\|_2}{\sqrt{n|y|}}.$$

Therefore, we have for all $\varepsilon \geq 4L \frac{\Lambda \max_{(x,y) \in S} \|\phi(x)\|_2}{\sqrt{2n|y|}}$

$$\log \left(\mathcal{N}_\infty(\varepsilon, \mathcal{G}, S) \right) \leq \frac{16 \max_{(x,y) \in S} \|\phi(x)\|_2^2 \Lambda^2 L^2}{\varepsilon^2} \log \left(\frac{2en|y|\hat{B}(S)L}{\varepsilon} \right).$$

Plugging this covering number in in our Rademacher bound with $\varepsilon \geq 4L \frac{\Lambda \max_{(x,y) \in S} \|\phi(x)\|_2}{\sqrt{2(|\mathcal{R}^X|+1)|y|}}$ and taking $S = \mathcal{R}^X \cup \{(X, Y)\}$ we get

$$\begin{aligned} \mathfrak{R}_{\mathcal{R}^X}(\mathcal{G}(X, Y)) &\leq \inf_{\varepsilon \in [0, \mathcal{G}_{\max}(X, Y)/2]} \left(4\varepsilon + \frac{12}{\sqrt{|\mathcal{R}^X|}} \int_\varepsilon^{\mathcal{G}_{\max}(X, Y)/2} \sqrt{\log \mathcal{N}_\infty(\nu/2, \mathcal{G}, \mathcal{R}^X \cup \{(X, Y)\})} d\nu \right) \\ &\leq \frac{16 \max_{(x,y) \in \mathcal{R}^X \cup \{(X, Y)\}} \|\phi(x)\|_2 \Lambda L}{\sqrt{2(|\mathcal{R}^X| + 1)|y|}} + \frac{12 \times 16 \max_{(x,y) \in \mathcal{R}^X \cup \{(X, Y)\}} \|\phi(x)\| \Lambda L}{\sqrt{|\mathcal{R}^X|}} \times \\ &\quad \times \int_{\frac{4L\Lambda \max_{(x,y) \in \mathcal{R}^X \cup \{(X, Y)\}} \|\phi(x)\|_2}{\sqrt{2(|\mathcal{R}^X|+1)|y|}}^{\mathcal{G}_{\max}(X, Y)/2} \frac{1}{\nu} \sqrt{\log \left(\frac{4e(|\mathcal{R}^X|+1)|y|\hat{B}(\mathcal{R}^X \cup \{(X, Y)\})L}{\nu} \right)} d\nu \\ &\leq \frac{16 \max_{(x,y) \in \mathcal{R}^X \cup \{(X, Y)\}} \|\phi(x)\|_2 \Lambda L}{\sqrt{2(|\mathcal{R}^X| + 1)|y|}} + \frac{8 \times 16 \max_{(x,y) \in \mathcal{R}^X \cup \{(X, Y)\}} \|\phi(x)\| \Lambda L}{\sqrt{|\mathcal{R}^X|}} \times \\ &\quad \times \left(\log \left(\frac{4\sqrt{2}eL\hat{B}(\mathcal{R}^X \cup \{(X, Y)\})(|\mathcal{R}^X|+1)|y|\sqrt{(|\mathcal{R}^X|+1)|y|}}{4L\Lambda \max_{(x,y) \in \mathcal{R}^X \cup \{(X, Y)\}} \|\phi(x)\|_2} \right) \right)^{3/2} \\ &\leq \frac{16 \max_{(x,y) \in \mathcal{R}^X \cup \{(X, Y)\}} \|\phi(x)\|_2 \Lambda L}{\sqrt{2(|\mathcal{R}^X| + 1)|y|}} + \frac{8 \times 16 \max_{(x,y) \in \mathcal{R}^X \cup \{(X, Y)\}} \|\phi(x)\| \Lambda L}{\sqrt{|\mathcal{R}^X|}} \times \\ &\quad \times \left(\log \left(\sqrt{2}e((|\mathcal{R}^X| + 1)|y|)^{3/2} \right) \right)^{3/2} \end{aligned}$$

In the final inequality we use the fact that

$$\begin{aligned} \hat{B}(\mathcal{R}^X \cup \{(X, Y)\}) &\leq \max_{(x,y) \in \mathcal{R}^X \cup \{(X, Y)\}} \|\phi(x)\|_2 \sup_{W: \|W\|_{2,2} \leq \Lambda} \|W\|_{2,\infty} \\ &\leq \max_{(x,y) \in \mathcal{R}^X \cup \{(X, Y)\}} \|\phi(x)\|_2 \Lambda \end{aligned}$$

Therefore, the final bound on the Rademacher complexity can be given as

$$\mathfrak{R}_{\mathcal{R}^X} \leq O \left(L_\ell \|\mathcal{F}^{\text{loc}}\|_\infty \frac{\ln(|y||\mathcal{R}^X|)^{3/2}}{\sqrt{|\mathcal{R}^X|}} \right). \quad (26)$$

Example: $\mathcal{F}^{\text{loc}} \equiv L$ -layer fully connected deep neural network (DNN)(Bartlett et al., 2017): Following Bartlett et al. (2017), we consider a L -layer deep neural network (DNN) $f_A = \sigma_L(A^L \sigma_{L-1}(A^{L-1} \sigma_{L-2}(\dots A^1 x)))$ for $x \in \mathcal{X}$ where

⁴We need $\tilde{\mathcal{F}}^X$ only to state some theorems in Lei et al. (2019). We refer interested readers to Lei et al. (2019) for the details.

$\mathcal{A} = (A_1, A_2, \dots, A_L)$ is the sequence of weight matrices. The matrix $A^l \in \mathbb{R}^{d_{l-1} \times d_l}$ for $l = 1$ to L , with $d_L = |\mathcal{Y}|$, and $d_0 = d$ given $\mathcal{X} \subseteq \mathbb{R}^d$. Furthermore, $\sigma_l(\cdot) : \mathbb{R}^{d_l} \rightarrow \mathbb{R}^{d_l}$ denotes the non-linearity (including pooling and activation), σ_l -s are taken to be 1-Lipschitz, and $\sigma_l(0) = 0$. We assume that the A^l matrix is initialized at M^l , for each $l = 1$ to L . We consider the local function class

$$\mathcal{F}^{\text{loc}} = \{f_{\mathcal{A}} : \|A^l - M^l\|_{2,1} \leq b_l, \|A^l\|_{\sigma} \leq s_l, \forall l \leq L-1\}.$$

Furthermore, we have for any $f \in \mathcal{F}^{\text{loc}}$ and any $x \in \mathcal{X}$ the function $(f(x), y) \rightarrow \ell(\gamma_f(\cdot, \cdot))$ is $2L_\ell$ -Lipschitz. Therefore, for a fixed set S , we have from Theorem 3.3 in Bartlett et al. (2017) that the covering number of the $\mathcal{G} = \{\ell(\gamma_f(\cdot, \cdot)) : f_{\mathcal{A}} \in \mathcal{F}^{\text{loc}}\}$ is given as

$$\log(\mathcal{N}_2(\varepsilon, \mathcal{G}, S)) \leq \frac{4L_\ell^2 B^2 \ln(2d_{\max}^2)}{\varepsilon^2} \left(\prod_{l=1}^L s_l\right)^2 \left(\sum_{l=1}^L (b_l/s_l)^{2/3}\right)^{3/2} = \frac{R}{\varepsilon^2},$$

where $d_{\max} = \max_{l=1}^L d_l$, $\sqrt{\frac{1}{|S|} \sum_{x \in S} \|x\|_2^2} \leq B$, and

$$R = 4L_\ell^2 B^2 \ln(2d_{\max}^2) \left(\prod_{l=1}^L s_l\right)^2 \left(\sum_{l=1}^L (b_l/s_l)^{2/3}\right)^{3/2}.$$

Using a the covering number based bound on Rademacher complexity we obtain

$$\begin{aligned} \mathfrak{R}_{\mathcal{R}^X}(\mathcal{G}(X, Y)) &\leq \inf_{\epsilon \in [0, \mathcal{G}_{2, \max}(X, Y)/2]} \left(4\epsilon + \frac{12}{\sqrt{|\mathcal{R}^X|}} \int_{\epsilon}^{\mathcal{G}_{2, \max}(X, Y)/2} \sqrt{\log\left(\frac{4L_\ell B \prod_{l=1}^L s_l}{\nu}\right) \log\left(\mathcal{N}_2\left(\nu/2, \mathcal{G}, \mathcal{R}^X\right)\right)} d\nu\right) \\ &\leq \inf_{\epsilon \in [0, \mathcal{G}_{2, \max}(X, Y)/2]} \left(4\epsilon + \frac{12}{\sqrt{|\mathcal{R}^X|}} \int_{\epsilon}^{\mathcal{G}_{2, \max}(X, Y)/2} \sqrt{\log\left(\frac{4L_\ell B \prod_{l=1}^L s_l}{\nu}\right) \frac{R}{\nu^2}} d\nu\right) \\ &\leq \inf_{\epsilon \in [0, \mathcal{G}_{2, \max}(X, Y)/2]} \left(4\epsilon + \frac{8\sqrt{R}}{\sqrt{|\mathcal{R}^X|}} \log^{3/2}\left(\frac{4L_\ell B \prod_{l=1}^L s_l}{\epsilon}\right) - \frac{8\sqrt{R}}{\sqrt{|\mathcal{R}^X|}} \log^{3/2}\left(\frac{8L_\ell B \prod_{l=1}^L s_l}{\mathcal{G}_{2, \max}(X, Y)}\right)\right) \\ &\leq \left(\frac{4\mathcal{G}_{2, \max}(X, Y)}{\sqrt{|\mathcal{R}^X|}} + \frac{8\sqrt{R}}{\sqrt{|\mathcal{R}^X|}} \log^{3/2}\left(\frac{4L_\ell B \prod_{l=1}^L s_l \sqrt{|\mathcal{R}^X|}}{\mathcal{G}_{2, \max}(X, Y)}\right)\right) - \frac{8\sqrt{R}}{\sqrt{|\mathcal{R}^X|}} \log^{3/2}\left(\frac{8L_\ell B \prod_{l=1}^L s_l}{\mathcal{G}_{2, \max}(X, Y)}\right) \end{aligned}$$

B.4. Function approximation

The following proposition states that the expected loss between \mathcal{F}^{loc} classes can be bounded by the L_1 regression error between these two.

Proposition B.5. *Let \mathcal{F}_1 and \mathcal{F}_2 be two classes for scorer functions, and loss ℓ satisfy Assumption 3.2 with Lipschitz constant L_ℓ . Then for any $x \in \mathcal{X}$ the following holds*

$$\begin{aligned} \min_{f \in \mathcal{F}_1} \mathbb{E}_{D^{x,r}}[\ell(f(X), Y)] &\leq \min_{f \in \mathcal{F}_2} \mathbb{E}_{D^{x,r}}[\ell(f(X), Y)] \\ &\quad + 2L_\ell \max_{f' \in \mathcal{F}_2} \min_{f \in \mathcal{F}_1} \mathbb{E}_{D^{x,r}}[\max_s |f_s(X) - f'_s(X)|]. \end{aligned}$$

Proof. We compare the loss from two arbitrary (measurable w.r.t. $D^{X,r}$) functions f and \tilde{f} next, where we are trying to approximate \tilde{f} with f .

$$\begin{aligned} \mathbb{E}[\ell(f(X), Y)] &= \mathbb{E}[\ell(\gamma_f(X, Y))] \\ &= \mathbb{E}[\ell(\gamma_{\tilde{f}}(X, Y))] + \mathbb{E}[\ell(\gamma_f(X, Y)) - \ell(\gamma_{\tilde{f}}(X, Y))] \\ &\leq \mathbb{E}[\ell(\gamma_{\tilde{f}}(X, Y))] + L_\ell \mathbb{E}[|\gamma_f(X, Y) - \gamma_{\tilde{f}}(X, Y)|] \\ &= \mathbb{E}[\ell(\tilde{f}(X), Y)] + L_\ell \mathbb{E}[|f_Y(X) - \tilde{f}_Y(X) - \max_{s \neq Y} f_s(X) + \max_{s \neq Y} \tilde{f}_s(X)|] \\ &\leq \mathbb{E}[\ell(\tilde{f}(X), Y)] + 2L_\ell \mathbb{E}[\max_s |f_s(X) - \tilde{f}_s(X)|] \end{aligned}$$

Let $f^{X,*} := \arg \min_{f \in \mathcal{F}^X} \mathbb{E}[\ell(f(X), Y)]$. Replacing the $\tilde{f} = f^{X,*}$ and taking minima over \mathcal{F}^{loc} on both sides, we obtain.

$$\min_{f \in \mathcal{F}^{\text{loc}}} \mathbb{E}[\ell(f(X), Y)] \leq \min_{f \in \mathcal{F}^X} \mathbb{E}[\ell(f(X), Y)] + 2L_\ell \min_{f \in \mathcal{F}^{\text{loc}}} \mathbb{E}[\max_s |f_s(X) - f_s^{X,*}(X)|].$$

□

Applying the above result, we have $\varepsilon_{\text{loc}} = 2L_\ell \min_{f \in \mathcal{F}^{\text{loc}}} \mathbb{E}[\max_s |f_s(X) - f_s^{X,*}(X)|]$ in Eq. (6).

A similar argument establishes that $\varepsilon_X = 2L_\ell \min_{f \in \mathcal{F}^X} \mathbb{E}[\max_s |f_s(X) - f_s^*(X)|]$ in Eq. (4), where the function f^* is the population minimizer of the loss over distribution D^X among the function class $\mathcal{F}^{\text{global}}$.

B.5. Proof of Proposition 3.6

Under the weak density condition, for any $r > 0$ we have $\mathbb{P}_D[|\mathcal{R}^X| = 0] = 0$. Furthermore, for any $N \geq 1$, and $x \in \mathcal{X}$

$$\begin{aligned} \mathbb{P}_D[|\mathcal{R}^X| < N] &\leq \mathbb{P}_D\left[\sum_{i=1}^n \mathbb{1}(d(X_i, x) \leq r) < N\right] \\ &\leq \mathbb{P}_D\left[\sum_{i=1}^n \mathbb{1}(d(X_i, x) \leq \min\{r, \delta_{\text{wdc}} \rho_D(x)^{-1/d}\}) < N\right] \\ &\leq \exp(-2(p(x, r) - N/n)^2 n) \end{aligned}$$

Let $p(x, r) := \min\{c_{\text{wdc}}^d \rho_D(x) r^d, c_{\text{wdc}}^d \delta_{\text{wdc}}^d\}$, then $\mathbb{P}_D(d(X_i, x) \leq \min\{r, \delta_{\text{wdc}} \rho_D(x)^{-1/d}\}) \geq p(x, r)$. Using Chernoff bound we obtain the final inequality for the above definition of $p(x, r)$. It can be shown that choosing

$$N = n \left(\min\{c_{\text{wdc}}^d \rho_D(x) r^d, c_{\text{wdc}}^d \delta_{\text{wdc}}^d\} - \sqrt{\frac{\log(1/\delta)}{2n}} \right),$$

we obtain $\mathbb{P}_D[|\mathcal{R}^X| < N] \leq \delta$ for any $\delta > 0$.

Recall, $\mathbb{P}_D[\rho_D(X) < f_\rho(\delta)] \leq \delta$ for any $\delta > 0$. Let

$$N(r, \delta) \geq n \left(\min\{c_{\text{wdc}}^d f_\rho(\delta) r^d, c_{\text{wdc}}^d \delta_{\text{wdc}}^d\} - \sqrt{\frac{\log(2/\delta)}{2n}} \right).$$

Then, we have

$$\begin{aligned} \mathbb{P}_D[|\mathcal{R}^X| < N(r, \delta)] &\leq \mathbb{P}_D[|\mathcal{R}^X| < N(r, \delta) | \rho_D(X) \geq f_\rho(\delta/2)] + \mathbb{P}_D[\rho_D(X) < f_\rho(\delta/2)] \\ &\leq \delta/2 + \delta/2 = \delta \end{aligned}$$

For the first term in the final inequality, we use the fact that for all $x \in \mathcal{X}$ such that $\rho_D(x) \geq \rho_D^{-1}(\delta/2)$, we have $\mathbb{P}_D[|\mathcal{R}^X| < N(r, \delta)] \leq \delta/2$. For the second term in the final inequality, we just use the definition of $f_\rho(\delta)$.

B.6. Computation of the function $f_\rho(\delta)$ in Proposition 3.6

For non-degenerate multi-dimensional Gaussian distributions we have

$$\rho_{N(\mu, \Sigma)}(x) = (2\pi)^{-d/2} |\Sigma|^{-1/2} \exp(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu))$$

Therefore, the level sets are given as

$$\begin{aligned} \mathbb{P}_{N(\mu, \Sigma)}[x : \rho_{N(\mu, \Sigma)}(x) \geq (2\pi)^{-d/2} |\Sigma|^{-1/2} \gamma] &= \int_{x: (x-\mu)^T \Sigma^{-1}(x-\mu) \leq 2 \ln(1/\gamma)} (2\pi)^{-d/2} |\Sigma|^{-1/2} \exp(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)) dx \\ &= \int_{x: (x-\mu)^T \Sigma^{-1}(x-\mu) \leq 2 \ln(1/\gamma)} (2\pi)^{-d/2} |\Sigma|^{-1/2} \int_{t=\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}^{\infty} \exp(-t) dt dx \end{aligned}$$

$$\begin{aligned}
 &= (2\pi)^{-d/2} |\Sigma|^{-1/2} \int_{t=0}^{\infty} \underbrace{\left(\int_{x: (x-\mu)^T \Sigma^{-1} (x-\mu) \leq \min\{2t, 2\ln(1/\gamma)\}} dx \right)}_{\text{volume of ellipsoid}} \exp(-t) dt \\
 &= (2\pi)^{-d/2} |\Sigma|^{-1/2} |\Sigma|^{1/2} \frac{\pi^{d/2}}{\Gamma(d/2+1)} \left(\int_{t=\ln(1/\gamma)}^{\infty} (2\ln(1/\gamma))^{d/2} \exp(-t) dt + \int_{t=0}^{\ln(1/\gamma)} (2t)^{d/2} \exp(-t) dt \right) \\
 &= \frac{1}{\Gamma(d/2+1)} \left(\Gamma(d/2+1) - \int_{\ln(1/\gamma)}^{\infty} (t^{d/2} - \ln(1/\gamma)^{d/2}) \exp(-t) dt \right) \\
 &= 1 - \frac{1}{\Gamma(d/2+1)} \int_0^{\infty} ((t' + \ln(1/\gamma))^{d/2} - \ln(1/\gamma)^{d/2}) \exp(-(t' + \ln(1/\gamma))) dt' \\
 &\geq 1 - \frac{\gamma}{\Gamma(d/2+1)} \int_0^{\infty} (q(t'/q)^{d/2} + r(\ln(1/\gamma)/r)^{d/2} - \ln(1/\gamma)^{d/2}) \exp(-t') dt' \quad [(q\alpha + r\beta)^p \leq q\alpha^p + r\beta^p : q + r = 1] \\
 &\geq 1 - \left(\ln(1/\gamma)^{d/2-1} \gamma + \frac{((1-\ln(1/\gamma))^{-d/2+1}-1)}{\Gamma(d/2+1)} \gamma \ln(1/\gamma)^{d/2} \right) \\
 &\geq 1 - \left(\ln(1/\gamma)^{d/2-1} \gamma + \frac{((1-\ln(1/\gamma))^{-d/2+1}-1)}{\Gamma(d/2+1)} \gamma \ln(1/\gamma)^{d/2} \right) \geq 1 - 2.45\gamma \ln(1/\gamma)^{d/2}, \forall \gamma \leq 1/2
 \end{aligned}$$

We now extend the results to mixture of distributions. It is easily shown below that if each of the mixture component $k \leq K$ satisfies $\mathbb{P}_{D_k} [x : \rho_{D_k}(x) \leq \gamma] \leq c\gamma \ln(1/\gamma)^{d/2}$ then

$$\begin{aligned}
 \mathbb{P}_{D_{\text{mix}}} [x : \rho_{D_{\text{mix}}}(x) \leq \gamma] &= \sum_k w_k \mathbb{P}_{D_k} [x : \sum_l w_l \rho_{D_l}(x) \leq \gamma] \leq \sum_k w_k \mathbb{P}_{D_k} [x : \rho_{D_k}(x) \leq \gamma w_k^{-1}] \\
 &\leq \sum_k w_k c\gamma w_k^{-1} \ln(w_k/\gamma)^{d/2} \leq cK\gamma \ln(1/\gamma)^{d/2}.
 \end{aligned}$$

C. Comparison of risk bounds

We now compare the risk bounds of the proposed explicit local ERM (which we think of as proxy towards understanding the implicit local learning happening in Retrieval augmented models) between different parametric, and non-parametric methods.

C.1. Sobolev spaces (Yarotsky, 2017)

In the following paragraph we briefly describe the setting of (Yarotsky, 2017) for the most part borrowing the notations from the authors. The authors study the approximation of functions $f : \mathbb{R}^d \rightarrow \mathbb{R}$, with Relu networks with the metrics $\max_{x \in [0,1]^d} |f(x) - \tilde{f}(x)|$ for some approximation \tilde{f} . They consider the Sobolev spaces $\mathcal{W}^{k,\infty}([0,1]^d)$ for $n = 1, 2, \dots$. For a function in Sobolev spaces $\mathcal{W}^{k,\infty}([0,1]^d)$ the weak derivatives upto order k are bounded in L_∞ norm. In particular, we define the norm in $\mathcal{W}^{k,\infty}([0,1]^d)$ as .

$$\|f\|_{\mathcal{W}^{k,\infty}([0,1]^d)} = \max_{|\mathbf{k}| \leq k} \text{ess sup}_{x \in [0,1]^d} \|D^{\mathbf{k}} f(x)\|_\infty,$$

where $\mathbf{k} \in \{0, 2, \dots, k\}^d$ is the multi-index of the weak derivative $D^{\mathbf{k}}$, and $|\mathbf{k}| = \sum_{k_i=1}^d k_i$.⁵ The function class

$$F_{k,d} = \{f \in \mathcal{W}^{k,\infty}([0,1]^d) : \|f\|_{\mathcal{W}^{k,\infty}([0,1]^d)} \leq 1\}.$$

From Theorem 4 we know that to approximate $F_{k,d}$ within accuracy $\epsilon \in (0, 1/2)$ we require $\Omega(\epsilon^{-d/2k})$ weights in general, and with a depth $O(\ln^p(1/\epsilon))$ network, for any $p \geq 0$, we require $\Omega(\epsilon^{-d/k} \ln^{-(2p+1)}(1/\epsilon))$ weights.

A standard bound of Taylor series ensures that a degree $(k-1)$ Taylor polynomial will approximate the function class $F_{k,d}$ for any x' in the L_2 -radius r of x (hence L_∞ -radius r as L_2 norm upper bounds L_∞ norm) with accuracy $d^k r^k / k!$, i.e. for any $f \in F_{k,d}$ there exists $\tilde{f}(x') \in \mathcal{P}(k)$

$$\max_{x': \|x'-x\|_2 \leq r} |f(x') - \tilde{f}(x')| \leq \frac{d^k r^k}{k!} \|f\|_{\mathcal{W}^{k,\infty}([0,1]^d)} \leq \frac{d^k r^k}{k!}.$$

⁵Note we adopt bold symbol only for multi-indices, whereas vectors in \mathbb{R}^d are denoted without bold symbol.

In particular, $\tilde{f}(x')$ can be taken as the Taylor polynomial of degree at most k $\tilde{f}(x') = \sum_{|\mathbf{k}| \leq k-1} a_{\mathbf{k}}(x - x')^{\mathbf{k}}$ where $|a_{\mathbf{k}}| \leq 1$. Hence, we have $\sum_{|\mathbf{k}| \leq k-1} a_{\mathbf{k}} \leq \binom{d+k-1}{k-1}$.

Connecting back to our approximation with the approximation error ε_{loc} we have for $\mathcal{F}^{\text{global}} = F_{k,d}$ the degree $q_x = (k-1)$ for all $x \in \mathcal{X}$, $C(F_{k,d}, k-1) = \frac{d^k r^k}{k!}$, and $C'(F_{k,d}, k-1) = \binom{d+k-1}{k-1}$.

D. Proofs for Section 3.4

This section focuses on providing a proof of Proposition 3.8. It follows the proof technique of (Foster et al., 2019, Eq. (9)). Before presenting the proof of Proposition 3.8, we need to introduce a slight variation of the Rademacher complexity for data-dependent hypothesis set.

Let $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$. Let $\mathcal{R} = \{z_j^{\mathcal{R}}\}, \mathcal{T} = \{z_j^{\mathcal{T}}\} \in \mathcal{Z}^m$ be two m -sized samples and $\sigma \in \{+1, -1\}^m$ be a vector of independent Rademacher variables. Now define $\mathcal{R}_{\mathcal{T}, \sigma} = \{z_j^{\mathcal{R}_{\mathcal{T}, \sigma}}\} \in \mathcal{Z}^m$ such that

$$z_j^{\mathcal{R}_{\mathcal{T}, \sigma}} = \begin{cases} z_j^{\mathcal{R}}, & \text{if } \sigma_j = 1, \\ z_j^{\mathcal{T}}, & \text{if } \sigma_j = -1, \end{cases} \quad (27)$$

i.e., $\mathcal{R}_{\mathcal{T}, \sigma}$ is obtained by replacing i -th element of \mathcal{R} by i -th element of \mathcal{T} iff $\sigma_i = -1$. Let $\mathcal{U} \in \mathcal{Z}^{n-m}$ be an $m-n$ -sized sample; for $\mathcal{R} \in \mathcal{Z}^m$, $\mathcal{S}_{\mathcal{R}} = \mathcal{U} \cup \mathcal{R} \in \mathcal{Z}^n$. Note that, following this notation, we have $\mathcal{S}_{\mathcal{R}_{\mathcal{T}, \sigma}} = \mathcal{U} \cup \mathcal{R}_{\mathcal{T}, \sigma}$. For $\mathcal{S} \in \mathcal{Z}^n$, let $\mathcal{H}(\mathcal{S})$ be a data dependent function class (hypothesis set), which does not depend on the ordering of the elements in \mathcal{S} .

Definition D.1 (Rademacher complexity for data-dependent function class). *Let $\mathcal{H} = \{\mathcal{H}(\mathcal{S})\}_{\mathcal{S} \in \mathcal{Z}^n}$ be a family of data dependent function classes. Given $\mathcal{R} = \{z_j^{\mathcal{R}}\}_{j \in [m]}$, $\mathcal{T} = \{z_j^{\mathcal{T}}\}_{j \in [m]} \sim \mathcal{D}^m$ and $\mathcal{U} = \{z_{m+i}^{\mathcal{U}}\}_{i \in [n-m]}$, the empirical Rademacher complexity $\mathfrak{R}_{\mathcal{U}, \mathcal{R}, \mathcal{T}}^{\circ}(\mathcal{H})$ and Rademacher complexity $\mathfrak{R}_{\mathcal{U}, m}^{\circ}(\mathcal{H})$ are defined as follows.*

$$\begin{aligned} \mathfrak{R}_{\mathcal{U}, \mathcal{R}, \mathcal{T}}^{\circ}(\mathcal{H}) &= \frac{1}{m} \mathbb{E}_{\sigma} \left[\sup_{h \in \mathcal{H}(\mathcal{S}_{\mathcal{R}_{\mathcal{T}, \sigma}})} \sum_{i=1}^m \sigma_i h(z_i^{\mathcal{T}}) \right] \\ \mathfrak{R}_{\mathcal{U}}^{\circ}(\mathcal{H}) &= \frac{1}{m} \mathbb{E}_{\mathcal{R}, \mathcal{T} \sim \mathcal{D}^m} \left[\sup_{h \in \mathcal{H}(\mathcal{S}_{\mathcal{R}_{\mathcal{T}, \sigma}})} \sum_{i=1}^m \sigma_i h(z_i^{\mathcal{T}}) \right] \end{aligned} \quad (28)$$

D.1. Proof of Proposition 3.8

We are now ready to establish the proof of Proposition 3.8. As discussed above, we extend the proof technique of (Foster et al., 2019, Eq. (9)) to obtain this result. Our setting differs from that of (Foster et al., 2019) as the local ERM objective only depends on the retrieve samples \mathcal{R}^x while the function class of interest $\mathcal{F}_{\mathcal{S}} = \mathcal{F}_{\Phi_{\mathcal{S}}}$ in (16) depends on the entire training set \mathcal{S} via representation $\Phi_{\mathcal{S}}$. We suitably modify the proof techniques of (Foster et al., 2019) to handle this difference.

Let $|\mathcal{R}^x| := m$ and $\mathcal{U} = \mathcal{S} \setminus \mathcal{R}^x$. For $\mathcal{R}, \mathcal{T} \in \mathcal{Z}^m$, we define

$$\begin{aligned} \Xi(\mathcal{R}, \mathcal{T}) &= \sup_{f \in \mathcal{F}_{\Phi_{\mathcal{U} \cup \mathcal{R}}}} \left| \underbrace{\mathbb{E}_{(X', Y') \sim \mathcal{D}^{x,r}}[\ell(f(X'), Y')]}_{:= R_{\ell}(f; \mathcal{D}^{x,r})} - \frac{1}{m} \underbrace{\sum_{(x', y') \in \mathcal{T}} \ell(f(x'), y')}_{:= \widehat{R}_{\ell}(f; \mathcal{T})} \right| \\ &= \sup_{f \in \mathcal{F}_{\Phi_{\mathcal{U} \cup \mathcal{R}}}} |R_{\ell}(f; \mathcal{D}^{x,r}) - \widehat{R}_{\ell}(f; \mathcal{T})|. \end{aligned}$$

Note that we are interested in bounding

$$\Xi(\mathcal{R}^x, \mathcal{R}^x) = \sup_{f \in \mathcal{F}_{\Phi_{\mathcal{S}}}} \left| \underbrace{\mathbb{E}_{(X', Y') \sim \mathcal{D}^{x,r}}[\ell(f(X'), Y')]}_{R_{\ell}(f; \mathcal{D}^{x,r})} - \frac{1}{m} \underbrace{\sum_{(x', y') \in \mathcal{T}} \ell(f(x'), y')}_{\widehat{R}_{\ell}(f; \mathcal{R}^x) = \widehat{R}_{\ell}^x(f)} \right|$$

where we have used the fact that $\mathcal{U} \cup \mathcal{R}^x = \mathcal{S}$. Towards this, we first establish that $\Xi(\mathcal{R}, \mathcal{R})$ satisfies the $(\frac{M_{\ell}}{m} + 2\Delta LL_{\ell,1})$ -bounded difference property, i.e., for $\mathcal{R}, \mathcal{R}' \in \mathcal{Z}^m$ that only differ in one element, we have

$$\Xi(\mathcal{R}, \mathcal{R}) - \Xi(\mathcal{R}', \mathcal{R}') \leq \frac{M_{\ell}}{m} + 2\Delta LL_{\ell,1}. \quad (29)$$

Note that

$$\Xi(\mathcal{R}, \mathcal{R}) - \Xi(\mathcal{R}', \mathcal{R}') \leq \underbrace{\Xi(\mathcal{R}, \mathcal{R}) - \Xi(\mathcal{R}, \mathcal{R}')}_{\text{I}} + \underbrace{\Xi(\mathcal{R}, \mathcal{R}') - \Xi(\mathcal{R}', \mathcal{R}')}_{\text{II}}. \quad (30)$$

Now, we will separately bound the two terms in the RHS. Let $\check{z} = (\check{x}, \check{y}) \in \mathcal{R} \setminus \mathcal{R}'$ and $\check{z}' = (\check{x}', \check{y}') \in \mathcal{R}' \setminus \mathcal{R}$. Thus, we have the following bound on the first term.

$$\begin{aligned} \text{I} &= \Xi(\mathcal{R}, \mathcal{R}) - \Xi(\mathcal{R}, \mathcal{R}') \\ &= \sup_{f \in \mathcal{F}_{\Phi_{\mathcal{U} \cup \mathcal{R}}}} |R_\ell(f; \mathcal{D}^{x,r}) - \widehat{R}_\ell(f; \mathcal{R})| - \sup_{f \in \mathcal{F}_{\Phi_{\mathcal{U} \cup \mathcal{R}}}} |R_\ell(f; \mathcal{D}^{x,r}) - \widehat{R}_\ell(f; \mathcal{R}')| \\ &\leq \sup_{f \in \mathcal{F}_{\Phi_{\mathcal{U} \cup \mathcal{R}}}} \left| |R_\ell(f; \mathcal{D}^{x,r}) - \widehat{R}_\ell(f; \mathcal{R})| - |R_\ell(f; \mathcal{D}^{x,r}) - \widehat{R}_\ell(f; \mathcal{R}')| \right| \\ &\leq \sup_{f \in \mathcal{F}_{\Phi_{\mathcal{U} \cup \mathcal{R}}}} [R_\ell(f; \mathcal{D}^{x,r}) - \widehat{R}_\ell(f; \mathcal{R}) - R_\ell(f; \mathcal{D}^{x,r}) + \widehat{R}_\ell(f; \mathcal{R}')] \\ &= \sup_{f \in \mathcal{F}_{\Phi_{\mathcal{U} \cup \mathcal{R}}}} |\widehat{R}_\ell(f; \mathcal{R}') - \widehat{R}_\ell(f; \mathcal{R})| \\ &= \sup_{f \in \mathcal{F}_{\Phi_{\mathcal{U} \cup \mathcal{R}}}} \frac{1}{m} |\ell(f(\check{x}'), \check{y}') - \ell(f(\check{x}), \check{y})| \leq \frac{M_\ell}{m}, \end{aligned} \quad (31)$$

where the last inequality follows from our boundedness assumption for the loss function ℓ .

Now we move to term II. Towards this, note that, it follows from the definition of supremum that, for any $\epsilon > 0$, there exists $\tilde{f} \in \mathcal{F}_{\Phi_{\mathcal{U} \cup \mathcal{R}'}}$ such that

$$\sup_{f \in \mathcal{F}_{\Phi_{\mathcal{U} \cup \mathcal{R}'}}} |R_\ell(f; \mathcal{D}^{x,r}) - \widehat{R}_\ell(f; \mathcal{R}')| - \epsilon \leq |R_\ell(\tilde{f}; \mathcal{D}^{x,r}) - \widehat{R}_\ell(\tilde{f}; \mathcal{R}')| \quad (32)$$

Let $\tilde{f} = \tilde{g} \circ \Phi_{\mathcal{U} \cup \mathcal{R}} \in \mathcal{F}_{\Phi_{\mathcal{U} \cup \mathcal{R}'}}$ and $\tilde{f}' = \tilde{g} \circ \Phi_{\mathcal{U} \cup \mathcal{R}'} \in \mathcal{F}_{\Phi_{\mathcal{U} \cup \mathcal{R}'}}$. Note that, for any $(x, y) \in \mathcal{Z}$,

$$\begin{aligned} |\ell(\tilde{f}(x), y) - \ell(\tilde{f}'(x), y)| &= |\ell(\tilde{g} \circ \Phi_{\mathcal{U} \cup \mathcal{R}}(x), y) - \ell(\tilde{g} \circ \Phi_{\mathcal{U} \cup \mathcal{R}'}(x), y)| \\ &\stackrel{(i)}{\leq} L_{\ell,1} \|\tilde{g} \circ \Phi_{\mathcal{U} \cup \mathcal{R}}(x) - \tilde{g} \circ \Phi_{\mathcal{U} \cup \mathcal{R}'}(x)\|_\infty \\ &\leq L_{\ell,1} \|\tilde{g} \circ \Phi_{\mathcal{U} \cup \mathcal{R}}(x) - \tilde{g} \circ \Phi_{\mathcal{U} \cup \mathcal{R}'}(x)\|_2 \\ &\stackrel{(ii)}{\leq} L_{\ell,1} L \|\Phi_{\mathcal{U} \cup \mathcal{R}}(x) - \Phi_{\mathcal{U} \cup \mathcal{R}'}(x)\|_2 \\ &\stackrel{(iii)}{\leq} L_{\ell,1} L \Delta, \end{aligned} \quad (33)$$

where we use $L_{\ell,1}$ -Lipschitzness of ℓ w.r.t. $\|\cdot\|_\infty$ norm, L -Lipschitzness of g , and Δ -sensitivity of the representation Φ in (i), (ii), and (iii), respectively.

Now, we have

$$\begin{aligned} \text{II} &= \Xi(\mathcal{R}, \mathcal{R}') - \Xi(\mathcal{R}', \mathcal{R}') \\ &= \sup_{f \in \mathcal{F}_{\Phi_{\mathcal{U} \cup \mathcal{R}}}} |R_\ell(f; \mathcal{D}^{x,r}) - \widehat{R}_\ell(f; \mathcal{R}')| - \sup_{f \in \mathcal{F}_{\Phi_{\mathcal{U} \cup \mathcal{R}'}}} |R_\ell(f; \mathcal{D}^{x,r}) - \widehat{R}_\ell(f; \mathcal{R}')| \\ &\stackrel{(i)}{\leq} |R_\ell(\tilde{f}; \mathcal{D}^{x,r}) - \widehat{R}_\ell(\tilde{f}; \mathcal{R}')| + \epsilon - \sup_{f \in \mathcal{F}_{\Phi_{\mathcal{U} \cup \mathcal{R}'}}} |R_\ell(f; \mathcal{D}^{x,r}) - \widehat{R}_\ell(f; \mathcal{R}')| \\ &\leq |R_\ell(\tilde{f}; \mathcal{D}^{x,r}) - \widehat{R}_\ell(\tilde{f}; \mathcal{R}')| + \epsilon - |R_\ell(\tilde{f}'; \mathcal{D}^{x,r}) - \widehat{R}_\ell(\tilde{f}'; \mathcal{R}')| \\ &= \left| [R_\ell(\tilde{f}; \mathcal{D}^{x,r}) - R_\ell(\tilde{f}'; \mathcal{D}^{x,r})] - [\widehat{R}_\ell(\tilde{f}; \mathcal{R}') - \widehat{R}_\ell(\tilde{f}'; \mathcal{R}')] \right| + \epsilon \\ &\leq |R_\ell(\tilde{f}; \mathcal{D}^{x,r}) - R_\ell(\tilde{f}'; \mathcal{D}^{x,r})| + |\widehat{R}_\ell(\tilde{f}; \mathcal{R}') - \widehat{R}_\ell(\tilde{f}'; \mathcal{R}')| + \epsilon \\ &\stackrel{(ii)}{\leq} 2L_{\ell,1} L \Delta + \epsilon, \end{aligned} \quad (34)$$

where (i) and (ii) follow from (32) and (33), respectively. Now, since ϵ in (32) can be chosen arbitrarily small, it follows from (30), (31), and (34) that

$$\Xi(\mathcal{R}, \mathcal{R}) - \Xi(\mathcal{R}', \mathcal{R}') \leq \frac{M_\ell}{m} + 2\Delta LL_{\ell,1},$$

i.e., $\Xi(\mathcal{R}, \mathcal{R})$ indeed satisfies the $(\frac{M_\ell}{m} + 2\Delta LL_{\ell,1})$ -bounded difference property. Now, it follows from the McDiarmid's inequality that, for $\delta > 0$, we have with probability at least $1 - \delta$:

$$\Xi(\mathcal{R}^x, \mathcal{R}^x) \leq \mathbb{E}[\Xi(\mathcal{R}^x, \mathcal{R}^x)] + (M_\ell + 2\Delta LL_{\ell,1}m) \sqrt{\frac{\log(1/\delta)}{2m}}$$

or

$$\begin{aligned} \sup_{f \in \mathcal{F}_{\Phi_S}} |R_\ell(f; D^{x,r}) - \widehat{R}_\ell^x(f)| &\leq \mathbb{E}_{\mathcal{R}^x} \left[\sup_{f \in \mathcal{F}_{\Phi_S}} |R_\ell(f; D^{x,r}) - \widehat{R}_\ell^x(f)| \right] + \\ &\quad (M_\ell + 2\Delta LL_{\ell,1}m) \sqrt{\frac{\log(1/\delta)}{2m}}. \end{aligned} \quad (35)$$

Now, first statement of Proposition 3.8 follows from (35) and the fact that $m = |\mathcal{R}^x|$.

It follows from the proof steps in Foster et al. (2019, Section E.1) that

$$\mathbb{E}_{\mathcal{R}^x} \left[\sup_{f \in \mathcal{F}_{\Phi_S = \mathcal{U} \cup \mathcal{R}^x}} |R_\ell(f; D^{x,r}) - \widehat{R}_\ell^x(f)| \right] \leq 2\mathfrak{R}_{\mathcal{U}}^\diamond(\ell \circ \mathcal{F}), \quad (36)$$

where $\mathcal{F} = \{\mathcal{F}_{\Phi_{\mathcal{U} \cup \mathcal{R}}}\}_{\mathcal{R} \in \mathcal{Z}^m}$ and $\mathfrak{R}_{\mathcal{U}}^\diamond$ is defined in (28). This completes the proof of Proposition 3.8. \square

E. Classification in extended feature space: A kernel-based approach

As introduced in Sec. 2.3, our objective is to learn a function $f : \mathcal{X} \times (\mathcal{X} \times \mathcal{Y})^* \rightarrow \mathbb{R}^{|\mathcal{Y}|}$. For a given instance x , such a function can leverage its neighboring set $\mathcal{R}^x \in (\mathcal{X} \times \mathcal{Y})^*$ to improve the prediction on x . In this work, we restrict ourselves to a sub-family of such retrieval-based methods that first map $\mathcal{R}^x \sim D^{x,r}$ to $\widehat{D}^{x,r}$ — an empirical estimate of the local distribution $D^{x,r}$, which is subsequently utilized to make a prediction for x . In particular, the scorers of interest are of the form:

$$(x, \mathcal{R}^x) \mapsto f(x, \widehat{D}^{x,r}) = (f_1(x, \widehat{D}^{x,r}), \dots, f_{|\mathcal{Y}|}(x, \widehat{D}^{x,r})) \in \mathbb{R}^{|\mathcal{Y}|}, \quad (37)$$

where $f_y(x, \widehat{D}^{x,r})$ denotes the score assigned to the y -th class. Thus, assuming that $\Delta_{\mathcal{X} \times \mathcal{Y}}$ denotes the set of distribution over $\mathcal{X} \times \mathcal{Y}$, we restrict to a suitable function class in $\{f : \mathcal{X} \times \Delta_{\mathcal{X} \times \mathcal{Y}} \rightarrow \mathbb{R}^{|\mathcal{Y}|}\}$. Note that, given a surrogate loss $\ell : \mathbb{R}^{|\mathcal{Y}|} \times \mathcal{Y} \rightarrow \mathbb{R}$ and scorer f , the empirical risk $\widehat{R}_\ell^{\text{ex}}(f)$ and population risk $R_\ell^{\text{ex}}(f)$ take the following form:

$$\widehat{R}_\ell^{\text{ex}}(f) = \frac{1}{n} \sum_{i \in [n]} \ell(x_i, \widehat{D}^{x_i,r}) \quad \text{and} \quad R_\ell^{\text{ex}}(f) = \mathbb{E}_{(X,Y) \sim D} [\ell(f(X, D^{X,r}), Y)]. \quad (38)$$

Note that that the general framework for learning in the extended feature space $\widetilde{\mathcal{X}} := \mathcal{X} \times \Delta_{\mathcal{X} \times \mathcal{Y}}$ provides a very rich class of functions. In this paper, we focus on a specific form of learning methods in the extended feature space by using the kernel methods. The method as well as its analysis is obtained by adapting the work on utilizing kernel methods for domain generalization (Blanchard et al., 2011; Deshmukh et al., 2019).

E.1. Kernel-based classification

Before introducing a kernel method for the classification, we need to define a suitable kernel $k : \widetilde{\mathcal{X}} \times \widetilde{\mathcal{X}} \rightarrow \mathbb{R}$ on the extended feature space $\widetilde{\mathcal{X}} := \mathcal{X} \times \Delta_{\mathcal{X} \times \mathcal{Y}}$. Towards this, let $k_{\mathcal{Z}}$ be a kernel over $\mathcal{Z} := \mathcal{X} \times \mathcal{Y}$. Assuming that $H_{k_{\mathcal{Z}}}$ is the reproducing kernel Hilbert space (RKHS) associated with $k_{\mathcal{Z}}$, we can define a kernel mean embedding (Smola et al., 2007) $\Psi : \Delta_{\mathcal{Z}} \rightarrow H_{k_{\mathcal{Z}}}$ as follows:

$$\Psi(P) = \int_{\mathcal{Z}} k_{\mathcal{Z}}(z, \cdot) dP. \quad (39)$$

For an empirical distribution $\hat{D}^{x,r}$ defined by \mathcal{R}^x , kernel embedding in (39) takes the following form.

$$\Psi(\hat{D}^{x,r}) = \frac{1}{|\mathcal{R}^x|} \sum_{(x',y') \in \mathcal{R}^x} k_{\mathcal{Z}}((x',y'), \cdot). \quad (40)$$

Now, using a kernel $k_{\mathcal{X}}$ over \mathcal{X} and a kernel-like function κ over $\Psi(\Delta_{\mathcal{Z}})$, we define a desired kernel $k : \tilde{\mathcal{X}} \times \tilde{\mathcal{X}} \rightarrow \mathbb{R}$ as follows:

$$k(\tilde{X}_1, \tilde{X}_2) = k((X_1, D^{X_1,r}), (X_2, D^{X_2,r})) = k_{\mathcal{X}}(X_1, X_2) \cdot \kappa(\Psi(D^{X_1,r}), \Psi(D^{X_2,r})). \quad (41)$$

Let H_k be the RKHS corresponding to the kernel k in (41), and $\|\cdot\|_{H_k}$ be the norm associated with H_k . Equipped with the kernel in (41) and associated H_k , for $\lambda > 0$, we propose to learn a scorer $f = (f_1, \dots, f_{|y|}) \in H_k^{|y|} := H_k \times \dots \times H_k$ via the following regularized ERM problem.

$$\hat{f}^{\text{ex}} = \arg \min_{f \in H_k^{|y|}} \frac{1}{n} \sum_{i=1}^n \ell(f(\tilde{x}_i), y_i) + \lambda \cdot \Omega(f), \quad (42)$$

where $\tilde{x}_i = (x_i, \hat{D}^{x_i,r})$ and $\Omega(f) := \|f\|_{H_k^{|y|}}^2 := \sum_{y \in \mathcal{Y}} \|f_y\|_{H_k}^2$. It follows from the representer theorem that the solution of (42) takes the form $\hat{f}^{\text{ex}}(\cdot) = \sum_{i \in [n]} \alpha_i k((x_i, \hat{D}^{x_i,r}), \cdot)$. One can apply multiclass extensions of SVMs to learn the weights $\{\alpha_i\}$ (Deshmukh et al., 2019). Next, we focus on studying the generalization behavior of the scorer \hat{f}^{ex} recovered in (42).

E.2. Generalization bounds for kernel-based classification

Before presenting a generalization bound for kernel-based classification over the extended feature space $\tilde{\mathcal{X}}$, we state the three key assumptions that are utilized in our analysis.

Assumption E.1. *The loss function $\ell : \mathbb{R}^{|y|} \times \mathcal{Y}$ is $L_{\ell,1}$ -Lipschitz w.r.t. the first argument, i.e.,*

$$|\ell(s_1, y) - \ell(s_2, y)| \leq L_{\ell,1} \cdot \|s_1 - s_2\|_{\infty} \quad \forall s_1, s_2 \in \mathbb{R}^{|y|} \text{ and } y \in \mathcal{Y}. \quad (43)$$

Furthermore, assume that $\sup_{(x,y)} \ell(x, y) := M_{\ell} \leq \infty$.

Assumption E.2. *Kernels $k_{\mathcal{X}}$, $k_{\mathcal{Z}}$, and κ are bounded by $M_{k_{\mathcal{X}}}$, $M_{k_{\mathcal{Z}}}$, and M_{κ} , respectively.*

Assumption E.3. *Let $H_{k_{\mathcal{Z}}}$ and H_{κ} be the RKHS associated with $k_{\mathcal{Z}}$ and κ , respectively. Then, the canonical feature map $\varphi_{\kappa} : H_{k_{\mathcal{Z}}} \rightarrow H_{\kappa}$ is α -Hölder continuous with $\alpha \in (0, 1]$, i.e.,*

$$\|\varphi_{\kappa}(h_1) - \varphi_{\kappa}(h_2)\|_{H_{\kappa}} \leq L' \cdot \|h_1 - h_2\|_{H_{k_{\mathcal{Z}}}}^{\alpha} \quad \forall h_1, h_2 \in \{h \in H_{k_{\mathcal{Z}}} : \|h\|_{H_{k_{\mathcal{Z}}}} \leq M_{k_{\mathcal{Z}}}\} \quad (44)$$

The following result states our generalization bound for the kernel-based classification method described in Sec. E.1.

Theorem E.4. *Let $0 \leq \delta \leq 1$ and Assumptions E.1–E.3 hold. Furthermore, let $N(r, \delta)$ be as defined in (15). Then, for any $B > 0$, the following holds with probability at least $1 - 3\delta$*

$$\begin{aligned} \sup_{f \in \mathcal{F}_B^k} |\hat{R}_{\ell}^{\text{ex}}(f) - R_{\ell}^{\text{ex}}(f)| &\leq 32\sqrt{\log 2} L_{\ell,1} B M_{k_{\mathcal{X}}} M_{k_{\mathcal{Z}}} n^{-\frac{1}{2}} \left(1 + \log^{\frac{3}{2}} \sqrt{2n|y|}\right) \\ &+ L_{\ell,1} L' M_{k_{\mathcal{X}}} B \left(M_{k_{\mathcal{Z}}} \sqrt{\frac{2 \log(\frac{n}{\delta})}{N(r, \frac{\delta}{n})}} + M_{k_{\mathcal{Z}}} \sqrt{\frac{1}{N(r, \frac{\delta}{n})}} + \frac{4M_{k_{\mathcal{Z}}} \log(\frac{n}{\delta})}{3N(r, \frac{\delta}{n})} \right)^{\alpha} + M \sqrt{\frac{\log(\frac{1}{\delta})}{2n}}, \end{aligned}$$

where $\mathcal{F}_B^k = \{f = (f_1, \dots, f_{|y|}) \in H_k^{|y|} : \Omega(f) \leq B^2\}$ and $M := M_{\ell} + L_{\ell,1} B M_{k_{\mathcal{X}}} M_{k_{\mathcal{Z}}}$.

Before presenting the proof of Theorem E.4, we state two key results from the literature that are used in our analysis.

Proposition E.5 ((Steinwart & Christmann, 2008)). *Let (Ω, \mathcal{A}, P) be a probability space, H be a separable Hilbert space, and $M > 0$. Let $\eta_1, \dots, \eta_m : \Omega \rightarrow H$ be m independent H -valued random variables satisfying $\|\eta_j\|_{\infty} \leq M$, for all $j \in [m]$. The, for $\delta > 0$, the following holds with probability at least $1 - \delta$.*

$$\left\| \frac{1}{m} \sum_{j=1}^m (\eta_j - \mathbb{E}_P[\eta_j]) \right\|_H \leq M \sqrt{\frac{2 \log(1/\delta)}{m}} + M \sqrt{\frac{1}{m}} + \frac{4M \log(1/\delta)}{3m}. \quad (45)$$

Proposition E.6. (Deshmukh et al., 2019; Lei et al., 2019) Let $\tilde{\mathcal{Z}} = \tilde{\mathcal{X}} \times \mathcal{Y}$ be (extended) input and output space pair and $\tilde{\mathcal{S}} = \{\tilde{z}_1, \dots, \tilde{z}_n\}$. Let H_k be a RKHS defined on $\tilde{\mathcal{X}}$, with k being the associated kernel. Let

$$\mathcal{F}_B^k = \{(f_1, \dots, f_{|\mathcal{Y}|}) : f_y \in H_k \forall y \in \mathcal{Y} \text{ and } \left(\sum_{y \in \mathcal{Y}} \|f_y\|_{H_k}^p \right)^{1/p} \leq B\}$$

and $\ell : \mathbb{R}^{|\mathcal{Y}|} \times \mathcal{Y} \rightarrow \mathbb{R}$ be a Lipschitz function in its first argument, i.e.,

$$|\ell(s_1, y) - \ell(s_2, y)| \leq L_{\ell,1} \|s_1 - s_2\|_{\infty} \quad \forall s_1, s_2 \in \mathbb{R}^{|\mathcal{Y}|} \text{ and } y \in \mathcal{Y}.$$

Then the Rademacher complexity of the induced function class $\ell \circ \mathcal{F}_B^k := \{\ell \circ f : f \in \mathcal{F}_B^k\}$ satisfies

$$\begin{aligned} \mathfrak{R}_{\tilde{\mathcal{S}}}(\ell \circ \mathcal{F}_B^k) &:= \mathbb{E}_{\sigma_i} \left[\sup_{f \in \mathcal{F}_B^k} \frac{1}{n} \sum_{i \in [n]} \sigma_i \ell(f(\tilde{x}_i), y_i) \right] \\ &\leq 16L_{\ell,1} \sqrt{\log 2B} \sup_{\tilde{x} \in \tilde{\mathcal{X}}} \sqrt{k(\tilde{x}, \tilde{x})} n^{-\frac{1}{2}} |\mathcal{Y}|^{\frac{1}{2} - \frac{1}{\max\{2, p\}}} \left(1 + \log^{\frac{3}{2}} \sqrt{2n} |\mathcal{Y}| \right). \end{aligned} \quad (46)$$

Note that $\sigma = (\sigma_1, \dots, \sigma_n)$ denotes n i.i.d. Rademacher random variable.

Proof of Theorem E.4. Note that

$$\begin{aligned} \sup_{f \in \mathcal{F}_B^k} \left| \widehat{R}_{\ell}^{\text{ex}}(f) - R_{\ell}^{\text{ex}}(f) \right| &= \sup_{f \in \mathcal{F}_B^k} \left| \frac{1}{n} \sum_{i=1}^n \ell(f(x_i, \widehat{D}^{x_i, r}), y_i) - \mathbb{E}_{(X, Y) \sim D} [\ell(f(X, D^{X, r}), Y)] \right| \\ &\leq \underbrace{\sup_{f \in \mathcal{F}_B^k} \left| \frac{1}{n} \sum_{i=1}^n \ell(f(x_i, \widehat{D}^{x_i, r}), y_i) - \frac{1}{n} \sum_{i=1}^n \ell(f(x_i, D^{x_i, r}), y_i) \right|}_{\text{I}} + \\ &\quad \underbrace{\sup_{f \in \mathcal{F}_B^k} \left| \frac{1}{n} \sum_{i=1}^n \ell(f(x_i, D^{x_i, r}), y_i) - \mathbb{E}_{(X, Y) \sim D} [\ell(f(X, D^{X, r}), Y)] \right|}_{\text{II}} \end{aligned} \quad (47)$$

Bounding the term-I in (47). Note that

$$\begin{aligned} \text{I} &= \sup_{f \in \mathcal{F}_B^k} \left| \frac{1}{n} \sum_{i=1}^n \ell(f(x_i, \widehat{D}^{x_i, r}), y_i) - \frac{1}{n} \sum_{i=1}^n \ell(f(x_i, D^{x_i, r}), y_i) \right| \\ &\leq \frac{L_{\ell,1}}{n} \sum_{i \in [n]} \|f(x_i, \widehat{D}^{x_i, r}) - f(x_i, D^{x_i, r})\|_{\infty} \\ &\leq \frac{L_{\ell,1}}{n} \sum_{i \in [n]} \max_{y \in \mathcal{Y}} |f_y(x_i, \widehat{D}^{x_i, r}) - f_y(x_i, D^{x_i, r})| \\ &\leq L_{\ell,1} \cdot \max_{y \in \mathcal{Y}} \max_{i \in [n]} |f_y(x_i, \widehat{D}^{x_i, r}) - f_y(x_i, D^{x_i, r})| \end{aligned} \quad (48)$$

It follows from the reproducing property of the kernel k that, for any $y \in \mathcal{Y}$,

$$\begin{aligned} |f_y(x_i, \widehat{D}^{x_i, r}) - f_y(x_i, D^{x_i, r})| &= |\langle f_y, k((x_i, \widehat{D}^{x_i, r}), \cdot) - k((x_i, D^{x_i, r}), \cdot) \rangle| \\ &\leq \|f_y\|_{H_k} \cdot \|k((x_i, \widehat{D}^{x_i, r}), \cdot) - k((x_i, D^{x_i, r}), \cdot)\|_{H_k}. \end{aligned} \quad (49)$$

Now,

$$\begin{aligned} &\|k((x_i, \widehat{D}^{x_i, r}), \cdot) - k((x_i, D^{x_i, r}), \cdot)\|_{H_k} \\ &= \left(k((x_i, \widehat{D}^{x_i, r}), (x_i, \widehat{D}^{x_i, r})) + k((x_i, D^{x_i, r}), (x_i, D^{x_i, r})) - 2k((x_i, \widehat{D}^{x_i, r}), (x_i, D^{x_i, r})) \right)_{H_k}^{1/2} \end{aligned}$$

$$= \sqrt{k_{\mathcal{X}}(x_i, x_i)} \left(\kappa(\Psi(\widehat{D}^{x_i, r}), \Psi(\widehat{D}^{x_i, r})) + \kappa(\Psi(D^{x_i, r}), \Psi(D^{x_i, r})) - 2\kappa(\Psi(\widehat{D}^{x_i, r}), \Psi(D^{x_i, r})) \right)_{H_k}^{1/2}$$

$$= \sqrt{k_{\mathcal{X}}(x_i, x_i)} \|\kappa(\Psi(\widehat{D}^{x_i, r}), \cdot) - \kappa(\Psi(D^{x_i, r}), \cdot)\|_{H_k}$$

$$\leq M_{k_{\mathcal{X}}} \|\kappa(\Psi(\widehat{D}^{x_i, r}), \cdot) - \kappa(\Psi(D^{x_i, r}), \cdot)\|_{H_{\kappa}} \quad (50)$$

$$= M_{k_{\mathcal{X}}} \|\varphi_{\kappa}(\Psi(\widehat{D}^{x_i, r})) - \varphi_{\kappa}(\Psi(D^{x_i, r}))\|_{H_{\kappa}}$$

$$\leq L' M_{k_{\mathcal{X}}} \cdot \|\Psi(\widehat{D}^{x_i, r}) - \Psi(D^{x_i, r})\|_{H_{k_z}}^{\alpha} \quad (51)$$

By combining (49) and (50), we obtain that

$$|f_y(x_i, \widehat{D}^{x_i, r}) - f_y(x_i, D^{x_i, r})| \leq L' M_{k_{\mathcal{X}}} \cdot \|f_y\|_{H_k} \cdot \|\Psi(\widehat{D}^{x_i, r}) - \Psi(D^{x_i, r})\|_{H_{k_z}}^{\alpha}. \quad (52)$$

Now, Hoeffding's inequality in Hilbert spaces (cf. Proposition E.5) implies that, for $i \in [n]$, the following holds with probability at least $1 - \delta$.

$$\|\Psi(\widehat{D}^{x_i, r}) - \Psi(D^{x_i, r})\|_{H_{k_z}}^{\alpha} = \left\| \frac{1}{|\mathcal{R}^{x_i}|} \sum_{(x', y') \in \mathcal{R}^{x_i}} k_z((x', y'), \cdot) - \mathbb{E}_{D^{x_i, r}} [k_z((X', Y'), \cdot)] \right\|_{H_{k_z}}$$

$$\leq M_{k_z} \sqrt{\frac{2 \log(1/\delta)}{|\mathcal{R}^{x_i}|}} + M_{k_z} \sqrt{\frac{1}{|\mathcal{R}^{x_i}|}} + \frac{4M_{k_z} \log(1/\delta)}{3|\mathcal{R}^{x_i}|}. \quad (53)$$

It follows from (52) and (53) that, for each $i \in [n]$,

$$|f_y(x_i, \widehat{D}^{x_i, r}) - f_y(x_i, D^{x_i, r})|$$

$$\leq L' M_{k_{\mathcal{X}}} \cdot \|f_y\|_{H_k} \cdot \left(M_{k_z} \sqrt{\frac{2 \log(1/\delta)}{|\mathcal{R}^{x_i}|}} + M_{k_z} \sqrt{\frac{1}{|\mathcal{R}^{x_i}|}} + \frac{4M_{k_z} \log(1/\delta)}{3|\mathcal{R}^{x_i}|} \right)^{\alpha} \quad \forall y \in \mathcal{Y} \quad (54)$$

holds with probability at least $1 - \delta$. Next, taking union bound over $i \in [n]$ implies that the following holds for all $i \in [n]$ and $y \in \mathcal{Y}$ with probability at least $1 - \delta$.

$$|f_y(x_i, \widehat{D}^{x_i, r}) - f_y(x_i, D^{x_i, r})|$$

$$\leq L' M_{k_{\mathcal{X}}} \|f_y\|_{H_k} \left(M_{k_z} \sqrt{\frac{2 \log(n/\delta)}{|\mathcal{R}^{x_i}|}} + M_{k_z} \sqrt{\frac{1}{|\mathcal{R}^{x_i}|}} + \frac{4M_{k_z} \log(n/\delta)}{3|\mathcal{R}^{x_i}|} \right)^{\alpha}. \quad (55)$$

Recall that, for each $i \in [n]$, we have $|\mathcal{R}^{x_i}| \geq N(r, \delta)$ with probability at least $1 - \delta$ (cf. (15)). Using union bound, we have $|\mathcal{R}^{x_i}| \geq N(r, \delta/n)$, $\forall i \in [n]$, with probability at least $1 - \delta$. Thus, the following holds for all $i \in [n]$ and $y \in \mathcal{Y}$ with probability at least $1 - 2\delta$

$$|f_y(x_i, \widehat{D}^{x_i, r}) - f_y(x_i, D^{x_i, r})|$$

$$\leq L' M_{k_{\mathcal{X}}} \|f_y\|_{H_k} \left(M_{k_z} \sqrt{\frac{2 \log(n/\delta)}{N(r, \delta/n)}} + M_{k_z} \sqrt{\frac{1}{N(r, \delta/n)}} + \frac{4M_{k_z} \log(n/\delta)}{3N(r, \delta/n)} \right)^{\alpha}. \quad (56)$$

By using $\|f_y\|_{H_k} \leq B$ and combining (48) with (56), we obtain that

$$\text{I} \leq L_{\ell, 1} L' M_{k_{\mathcal{X}}} B \left(M_{k_z} \sqrt{\frac{2 \log(n/\delta)}{N(r, \delta/n)}} + M_{k_z} \sqrt{\frac{1}{N(r, \delta/n)}} + \frac{4M_{k_z} \log(n/\delta)}{3N(r, \delta/n)} \right)^{\alpha} \quad (57)$$

holds with probability at least $1 - 2\delta$.

Bounding the term-II in (47). Note that

$$\text{II} = \sup_{f \in \mathcal{F}_B^k} \left| \frac{1}{n} \sum_{i=1}^n \ell(f(x_i, D^{x_i, r}), y_i) - \mathbb{E}_{(X, Y) \sim D} [\ell(f(X, D^{X, r}), Y)] \right| \quad (58)$$

Using the Assumptions E.1 and E.2 and the fact that $f \in \mathcal{F}_B^k$, we can argue that

$$\begin{aligned}
 \ell(f(x, D^{x,r}), y) &= \ell(0, y) + |\ell(f(x, D^{x,r}), y) - \ell(0, y)| \\
 &\leq M_\ell + L_{\ell,1} \|f(x, D^{x,r})\|_\infty \\
 &\leq M_\ell + L_{\ell,1} \max_{y' \in \mathcal{Y}} |\langle f_{y'}, k((x, D^{x,r}), \cdot) \rangle| \\
 &\leq M_\ell + L_{\ell,1} \max_{y' \in \mathcal{Y}} \|f_{y'}\|_{H_k} M_k \\
 &\leq M_\ell + L_{\ell,1} R M_k \leq M_\ell + L_{\ell,1} R M_{k_x} M_\kappa := M
 \end{aligned}$$

Now, it follows from the Azuma-McDiarmid's inequality that the following holds with probability at least $1 - \delta$.

$$\begin{aligned}
 &\sup_{f \in \mathcal{F}_B^k} \left| \frac{1}{n} \sum_{i=1}^n \ell(f(x_i, D^{x_i,r}), y_i) - \mathbb{E}_{(X,Y) \sim \mathcal{D}} [\ell(f(X, D^{X,r}), Y)] \right| \\
 &\leq \mathbb{E} \left[\sup_{f \in \mathcal{F}_B^k} \left| \frac{1}{n} \sum_{i=1}^n \ell(f(x_i, D^{x_i,r}), y_i) - \mathbb{E}_{(X,Y) \sim \mathcal{D}} [\ell(f(X, D^{X,r}), Y)] \right| \right] \\
 &\quad + M \sqrt{\frac{\log(1/\delta)}{2n}}, \tag{59}
 \end{aligned}$$

Using the standard symmetrization procedure, we get that

$$\begin{aligned}
 &\mathbb{E} \left[\sup_{f \in \mathcal{F}_B^k} \left| \frac{1}{n} \sum_{i=1}^n \ell(f(x_i, D^{x_i,r}), y_i) - \mathbb{E}_{(X,Y) \sim \mathcal{D}} [\ell(f(X, D^{X,r}), Y)] \right| \right] \\
 &\leq \frac{2}{n} \cdot \mathbb{E}_{(X_i, Y_i) \sim \mathcal{D}} \mathbb{E}_{\sigma_i} \left[\sum_{i \in [n]} \sigma_i \ell(f(x_i, D^{x_i,r}), y_i) \right], \\
 &= 2 \bar{\mathfrak{R}}_{\mathfrak{S}}(\ell \circ \mathcal{F}_B^k)
 \end{aligned}$$

where $\sigma = (\sigma_1, \dots, \sigma_n)$ denotes n i.i.d. Rademacher random variables and $\bar{\mathfrak{R}}_{\mathfrak{S}}(\ell \circ \mathcal{F}_B^k)$ denote the Rademacher complexity of the function class

$$\ell \circ \mathcal{F}_B^k = \left\{ (x, y, D^{x,r}) \mapsto \ell(f(x, D^{x,r}), y) : f \in \mathcal{F}_B^k \right\}.$$

Now, using Proposition E.6 with $p = 2$ and Assumption E.2, we have

$$\begin{aligned}
 \bar{\mathfrak{R}}_{\mathfrak{S}}(\ell \circ \mathcal{F}_B^k) &\leq 16L_{\ell,1} \sqrt{\log 2B} \sup_{\tilde{x} \in \tilde{\mathcal{X}}} \sqrt{k(\tilde{x}, \tilde{x})} n^{-\frac{1}{2}} \left(1 + \log^{\frac{3}{2}} \sqrt{2n} |\mathcal{Y}| \right) \\
 &\leq 16L_{\ell,1} \sqrt{\log 2B} M_\kappa M_{k_x} n^{-\frac{1}{2}} \left(1 + \log^{\frac{3}{2}} \sqrt{2n} |\mathcal{Y}| \right) \tag{60}
 \end{aligned}$$

Now, by combining (58), (59), and (60), we obtain that with probability at least $1 - \delta$

$$\Pi \leq 32 \sqrt{\log 2} L_{\ell,1} B M_\kappa M_{k_x} n^{-\frac{1}{2}} \left(1 + \log^{\frac{3}{2}} \sqrt{2n} |\mathcal{Y}| \right) + M \sqrt{\frac{\log(1/\delta)}{2n}}. \tag{61}$$

Finally, combining (47), (57) and (61) completes the proof. \square

E.3. Empirical verification

We run an experiment to empirically verify the kernel based extended feature space-based approach. We design a kernel for extended feature space as in (41). In particular, we use Gaussian-like function for

$$\kappa(\Psi(D^{x_1,r}), \Psi(D^{x_2,r})) = \exp(-\|\Psi(D^{x_1,r}) - \Psi(D^{x_2,r})\|^2 / 2\sigma_\kappa^2).$$

To empirically estimate the distance between kernel mean embeddings of the two distributions $\|\Psi(\hat{D}^{x_1,r}) - \Psi(\hat{D}^{x_2,r})\|^2$ we follow Muandet et al. (2017); Li et al. (2015) as:

$$\begin{aligned} \|\Psi(\hat{D}^{x_1,r}) - \Psi(\hat{D}^{x_2,r})\|^2 &= \frac{1}{|\mathcal{R}^{x_1}|^2} \sum_{(x',y') \in \mathcal{R}^{x_1}} \sum_{(x'',y'') \in \mathcal{R}^{x_1}} k_{\mathcal{Z}}((x',y'),(x'',y'')) \\ &\quad + \frac{1}{|\mathcal{R}^{x_2}|^2} \sum_{(x',y') \in \mathcal{R}^{x_2}} \sum_{(x'',y'') \in \mathcal{R}^{x_2}} k_{\mathcal{Z}}((x',y'),(x'',y'')) \\ &\quad - \frac{2}{|\mathcal{R}^{x_1}||\mathcal{R}^{x_2}|} \sum_{(x',y') \in \mathcal{R}^{x_1}} \sum_{(x'',y'') \in \mathcal{R}^{x_2}} k_{\mathcal{Z}}((x',y'),(x'',y'')) \end{aligned}$$

We took $k_{\mathcal{Z}}((x',y'),(x'',y'')) = \exp(-\|x' - x''\|/2\sigma_x^2 - \lambda \mathbb{1}\{y' \neq y''\})$, which is basically like a normal L2 distance with labels concatenated as one-hot vectors. Also, $k_{\mathcal{X}}(x_1, x_2) = \exp(-\|x_1 - x_2\|^2/2\sigma_x^2)$ with normal L2 distance, with which we finally obtain the overall kernel for the extended feature space as:

$$k(\tilde{X}_1, \tilde{X}_2) = k((X_1, D^{X_1,r}), (X_2, D^{X_2,r})) = k_{\mathcal{X}}(X_1, X_2) \cdot \kappa(\Psi(D^{X_1,r}), \Psi(D^{X_2,r})).$$

For the synthetic dataset from Sec. 5, results are tabulated below:

Table 1. Accuracy of kernel classifier over extended feature space kernel as a function of number of retrieved neighbors used to form the extended feature space.

Neighbors	Kernel machine
2	0.776 ± 0.023
5	0.769 ± 0.021
10	0.777 ± 0.019
20	0.835 ± 0.021
50	0.819 ± 0.021
100	0.792 ± 0.022
200	0.585 ± 0.020

Note that, as the generalization bound suggested, that the model performance only improves up to a specific number of neighbors and starts degrading when further increasing the number of neighbors.

It’s also worth highlighting that a (4 layer) Transformer model that directly processes an instance along with the associated retrieved neighboring examples achieves a much higher performance of 0.898 with 10 neighbors. This is consistent with similar observations in the deep learning literature where kernel-based methods are often significantly outperformed by end-to-end neural networks (Bai & Lee, 2019; Chen et al., 2020; Samarin et al., 2020).

F. Additional details for experiments

F.1. Synthetic

Task and data. We consider the task of binary classification on mixtures using *synthetic data*: In particular, we assume $k = 100$ clusters in a $D = 10$ -dimensional space. Each cluster is specified by a mean parameter $\mu_i \in \mathbb{R}^D \sim \text{Uniform}(-10, 10)$ and a classification weight vector $w_i \in \mathbb{R}^D \sim \mathcal{N}(0, \mathbb{I})$ for $i = 1, 2, \dots, k$. We randomly generate a train set of $n = 10000$ points as follows: To generate a labeled example $(x_j, y_j), j \in [n]$: 1) select a cluster i uniformly at random, and 2) sample $x_j \sim \mathcal{N}(\mu_i, \mathbb{I})$ and its label $y_j = \text{sign}(w_i^T(x_j - \mu_i))$. Additionally, we also generate another set of points as test set using the same procedure.

Methods. As baseline, we consider models of various complexity, starting from simple linear classifier, to support vector machines with polynomial kernel (of degree 3) and with radial basis function (RBF) kernel, to a multi-layer perceptron (MLP) of two layers. For retrieval-based models, we consider each of the above method as the local model to fit on retrieved data points via local ERM framework (Sec. 3). Additionally, we also report simple kNN baseline. We compare all these methods using classification accuracy on the held out test set. We repeat all the experiments 10 times.

Observations. In Figure 4, we observe the tradeoff of varying the size of the retrieved set (as dictated by the neighborhood radius) on the performance of the proposed algorithms. We see that when the number of retrieved samples is small the local methods have lower accuracy, this is due to large generalization error. When the size of the retrieved sample space is high, the local methods fail to minimize the loss effectively due to the lack of model capacity. We see that this effect being more pronounced for simpler function classes such as linear classifier as compared to RBF or polynomial classifiers.

F.2. CIFAR-10

Task and data. We consider the task of binary classification on a *real image data* for object detection. In particular, we consider a subset of CIFAR-10 dataset where we only restrict to images from ‘‘Cat’’ and ‘‘Dog’’ classes. We randomly partition the data into a train set of $n = 10000$ points and remaining 2000 points for test. We do a 10-fold cross-validation.

Methods. We consider a subset of method from Appendix. F.1. In particular, we only consider a simple linear classifier and a multi-layer perceptron (MLP) of two layers. For retrieval-based models, we consider each of the above methods as the local model to fit on retrieved data points via local ERM framework (Sec. 3). The retrieval is done using L2 distance in the input space directly (no features is extracted). Additionally, we also report simple kNN baseline. We compare all these methods using classification accuracy on the held out test set. We repeat all the experiments 10 times.

Observations. Similar to Figure 4, Figure 5 exhibits a tradeoff, where varying the size of the retrieved set (as dictated by the neighborhood radius) impacts the performance of the proposed algorithms. We see when the number of retrieved samples is small the local methods have lower accuracy, this is due to large generalization error; and when the number of retrieved samples is large, simple local function class incurs a large approximation error.

F.3. ImageNet

Task and data. We consider the task of 1000-way image classification on ImageNet ILSVRC-12 dataset. We use the standard train-test set split, where we have of $n = 1281167$ points for training and 50000 points for test. Given large computational cost, we could only run each experiment once.

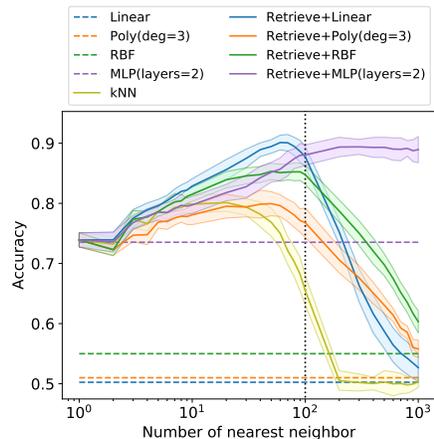


Figure 4. Performance of ERM and local ERM for various models on synthetic data.

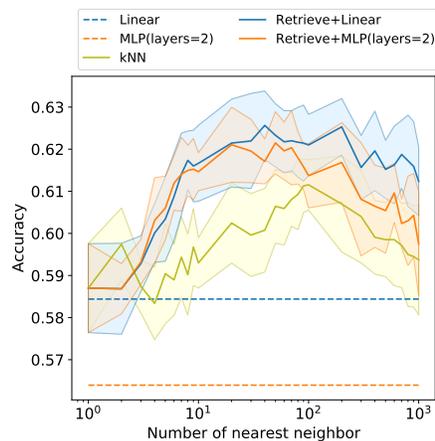


Figure 5. Performance of ERM and local ERM for various models on (binary) CIFAR-10.

Methods. We compare proposed Local ERM (Sec. 3) to state-of-the-art (SoTA) single model published for this task, which is from the most recent CVPR 2022 (Zhai et al., 2022). For the local parametric model we use a small MobileNetV3 architecture (Howard et al., 2019) with 4.01M parameters and 156 MFLOPs compute cost. Contrast this to SoTA model ViT-G/14 with 1.84B parameters and 938 GFLOPs compute cost. Following standard practice in literature, we use unsupervised learned features from ALIGN (Jia et al., 2021) to do image retrieval using L2 distance. For solving the local ERM, we fine-tune a MobileNetV3 model, which has been pretrained on ImageNet, on the retrieved set using Adam optimizer with a linear decay schedule. Additionally, we also report simple kNN baseline. We compare all these methods using classification accuracy on the held out test set.

Observations In Figure 6, we see that local ERM with a small MobileNet-V3 model is able to achieve the top-1 accuracy of 82.78 whereas a regularly trained MobileNet-V3 model achieves the top-1 accuracy of only 65.80. Also the result is very competitive with SoTA of 90.45 with a *much larger model*. Thus, the result suggest that the simple local ERM framework (analyzed in our work) is able to demonstrate the utility of retrieval-based models. In particular, it allows a realistic small sized model to attain very competitive numbers on the popular ImageNet benchmark. Furthermore, as pointed at end of Sec. 3.4, using global representation from ALIGN embeddings help simplest linear model to outperform MobileNet-V3 working directly on image input, thereby showcasing the benefits of endowing local ERM with global representation.

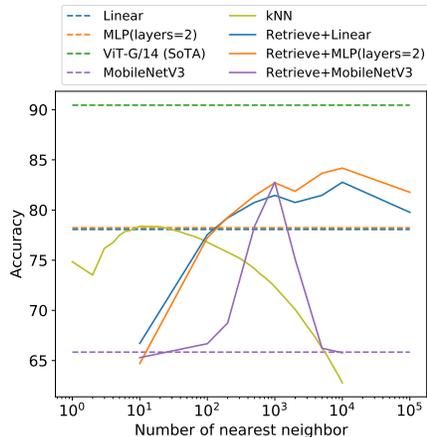


Figure 6. Performance of ERM and local ERM for various models on on ImageNet.