

Learning Infinite-Layer Networks. Beyond the Kernel Trick.

Amir, Globerson
gamir@cs.tau.ac.il

Roi, Livni
roi.livni@mail.huji.ac.il

June 17, 2016

Abstract

Infinite-Layer Networks (ILN) have recently been proposed as an architecture that mimics neural networks while enjoying some of the advantages of kernel methods. ILN are networks that integrate over infinitely many nodes within a single hidden layer. It has been demonstrated by several authors that the problem of learning ILN can be reduced to the kernel trick, implying that whenever a certain integral can be computed analytically they are efficiently learnable.

In this work we give an online algorithm for ILN, which avoids the *kernel trick assumption*. More generally and of independent interest, we show that kernel methods in general can be exploited even when the kernel cannot be efficiently computed but can only be estimated via sampling.

We provide a regret analysis for our algorithm, showing that it matches the sample complexity of methods which have access to kernel values. Thus, our method is the first to demonstrate that the kernel trick is not necessary as such, and random features suffice to obtain comparable performance.

1 Introduction

With the increasing success of highly non-convex and complex learning architectures such as neural networks, there is an increasing effort to further understand and explain the limits of training such hierarchical structures.

Recently there have been attempts to draw mathematical insight from kernel methods in order to better understand deep learning as well as come up with new computationally learnable architectures. One such line of work consists on learning *Infinite-Layer Networks* (ILN) [8, 16]. An infinite layer network can be thought of as a network with infinitely many nodes in a hidden layer. A target function in an ILN class will be of the form:

$$\mathbf{x} \rightarrow \int \psi(\mathbf{x}; \mathbf{w}) f(\mathbf{w}) d\mu(\mathbf{w}), \quad (1)$$

Here ψ is some function of the input \mathbf{x} and parameters \mathbf{w} , and $d\mu(\mathbf{w})$ is a prior over the parameter space. For example, $\psi(\mathbf{x}; \mathbf{w})$ can be a single sigmoidal neuron or a complete

convolutional network. The integral can be thought of as an infinite sum over all such possible networks, and $f(\mathbf{w})$ can be thought of as an infinite output weight vector to be trained.

A Standard 1-hidden layer network with a finite set of units can be obtained from the above formalism as follows. First, choose $\psi(\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{x} \cdot \mathbf{w})$ where σ is an activation function (e.g., sigmoid or relu). Next, set $d\mu(\mathbf{w})$ to be a discrete measure over a finite set $\mathbf{w}_1, \dots, \mathbf{w}_d$.¹ In this case, the integral results in a network with d hidden units, and the function f is the linear weights of the output layer. Namely:

$$\mathbf{x} \rightarrow \frac{1}{d} \sum_{i=1}^d f(\mathbf{w}_i) \cdot \sigma(\mathbf{x} \cdot \mathbf{w}_i).$$

The main challenge when training 1-hidden layer networks is of course to *find* the $\mathbf{w}_1, \dots, \mathbf{w}_d$ on which we wish to support our distribution. It is known [20], that due to hardness of learning intersection of halfspaces [19, 12], 1-hidden layer neural networks are computationally hard for a wide class of activation functions. Therefore, as the last example illustrates, the choice of μ is indeed crucial for performance.

For a fixed prior μ , the class of ILN functions is highly expressive, since f can be chosen to approximate any 1-hidden layer architecture to arbitrary precision (by setting f to delta functions around the weights of the network, as we did above for μ). However, this expressiveness comes at a cost. As argued in [16], ILNs will generalize well when there is a large mass of \mathbf{w} parameters that attain a small loss.

The key observation that makes certain ILN tractable to learn is that Eq. 1 is a linear functional in f . In that sense it is a linear classifier and enjoys the rich theory and algorithmic toolbox for such classifiers. In particular, one can use the fact that linear classifiers can be learned via the kernel trick. In other words, we can reduce learning ILN to the problem of computing the kernel function between two examples. Specifically the problem reduces to computing integrals of the following form:

$$k(\mathbf{x}_1, \mathbf{x}_2) = \int \psi(\mathbf{x}_1; \mathbf{w}) \cdot \psi(\mathbf{x}_2; \mathbf{w}) d\mu(\mathbf{w}) = \mathbb{E}_{\mathbf{w} \sim \mu} [\psi(\mathbf{x}_1; \mathbf{w}) \cdot \psi(\mathbf{x}_2; \mathbf{w})]. \quad (2)$$

In this work we extend this result to the case where no closed form kernel is available, and thus the kernel trick is not directly applicable. We thus turn our attention to the setting where features (i.e., \mathbf{w} vectors) can be randomly sampled. In this setting, our main result shows that for the square loss, we can efficiently learn the above class – and surprisingly with a comparable computational cost w.r.t standard kernel methods. Our approach can thus be applied to any random feature based method.

The observation we begin with is that sampling random neurons, i.e. sampling random \mathbf{w} , leads to an estimate of the kernel in Eq. 2. Thus, if for example, we ignore complexity issues and can sample infinitely many \mathbf{w} 's, it is not surprising that we can avoid the need for exact computation of the kernel.

¹In δ function notation $d\mu(\mathbf{w}) = \frac{1}{d} \sum_{i=1}^d \delta(\mathbf{w} - \mathbf{w}_i) d\mathbf{w}$

Our results provide a much stronger and practical result. Given T training samples, the lower bound on achievable accuracy is $O(1/\sqrt{T})$ (see [26]). We show that we can in fact achieve this rate, using $\tilde{O}(T^2)$ calls² to the random feature generator. For comparison, note that $O(T^2)$ is the size of the kernel matrix, and is thus likely to be the cost of any algorithm that uses an explicit kernel matrix, where one is available (whether a kernel regression algorithm can use less than $O(T^2)$ kernel entries is an open problem [5]. As we discuss later, our approach improves on previous random features based learning [9, 24] in terms of sample/computational complexity, and expressiveness.

2 Problem Setup

We consider learning algorithms that learn a mapping from input instances $\mathbf{x} \in \mathcal{X}$ to labels $y \in \mathcal{Y}$. We focus on the regression case where \mathcal{Y} is the interval $[-1, 1]$. Our starting point is a class of feature functions $\psi(\mathbf{w}; \mathbf{x}) : \Omega \times \mathcal{X} \rightarrow \mathbb{R}$, parametrized by vectors $\mathbf{w} \in \Omega$. The functions $\psi(\mathbf{w}; \mathbf{x})$ may contain highly complex non linearities, such as multi-layer networks consisting of convolution and pooling layers. Our only assumption on $\psi(\mathbf{w}; \mathbf{x})$ is that for all $\mathbf{w} \in \Omega$ and $\mathbf{x} \in \mathcal{X}$ it holds that $|\psi(\mathbf{w}; \mathbf{x})| < 1$.

Given a distribution μ on Ω , we denote by $L_2(\Omega, \mu)$ the class of square integrable functions over Ω .

$$L_2(\Omega, \mu) = \left\{ f : \int f^2(\mathbf{w}) d\mu(\mathbf{w}) < \infty \right\}.$$

We will use functions $f \in L_2(\Omega, \mu)$ as mixture weights over the class Ω , where each f naturally defines a new regression function from \mathbf{x} to \mathbb{R} as follows:

$$\mathbf{x} \rightarrow \int \psi(\mathbf{w}; \mathbf{x}) f(\mathbf{w}) d\mu(\mathbf{w}). \quad (3)$$

Our key algorithmic assumption is that the learner can efficiently sample random \mathbf{w} according to the distribution μ . Denote the time to generate one such sample by ρ .

In what follows it will be simpler to express the integrals as scalar products. Define the following scalar product on functions $f \in L_2(\Omega, \mu)$.

$$\langle f, g \rangle = \int f(\mathbf{w}) g(\mathbf{w}) d\mu(\mathbf{w}) \quad (4)$$

We denote the corresponding ℓ_2 norm by $\|f\| = \sqrt{\langle f, f \rangle}$. Also, given features \mathbf{x} denote by $\Phi(\mathbf{x})$ the function in $L_2(\Omega, \mu)$ given by $\Phi(\mathbf{x})[\mathbf{w}] = \psi(\mathbf{w}; \mathbf{x})$. Then, the regression functions we are considering are of the form $\mathbf{x} \rightarrow \langle f, \Phi(\mathbf{x}) \rangle$.

A subclass of norm bounded elements in $L_2(\Omega, \mu)$ induces a natural subclass of regression functions. Namely, we consider the following class:

$$\mathcal{H}_\mu^B = \{ \mathbf{x} \rightarrow \langle f, \Phi(\mathbf{x}) \rangle : \|f\| < B \}.$$

²We use \tilde{O} notation to suppress logarithmic factors

Our ultimate goal is to output a predictor $f \in L_2(\Omega, \mu)$ that is competitive, in terms of prediction, with the best target function in the class \mathcal{H}_μ^B .

We will consider an online setting, and use it to derive generalization bounds via standard online to batch conversion. In our setting, at each round a learner chooses a target function $f_t \in L_2(\Omega, \mu)$ and an adversary then reveals a sample \mathbf{x}_t and label y_t . The learner then incurs a loss of

$$\ell_t(f_t) = \frac{1}{2} (\langle f_t, \Phi(\mathbf{x}_t) \rangle - y_t)^2. \quad (5)$$

The objective of the learner is to minimize her T round regret w.r.t norm bounded elements in $L_2(\Omega, \mu)$. Namely:

$$\sum_{t=1}^T \ell_t(f_t) - \min_{f^* \in \mathcal{H}_\mu^B} \sum_{t=1}^T \ell_t(f^*). \quad (6)$$

In the statistical setting we assume that the sequence $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^T$ is generated IID according to some unknown distribution \mathbb{P} . We then define the expected loss of a predictor as

$$L(f) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}} \left[\frac{1}{2} (\langle f, \Phi(\mathbf{x}) \rangle - y)^2 \right]. \quad (7)$$

3 Main Results

Theorem 1 states our result for the online model. The corresponding result for the statistical setting is given in Corollary 1. We will elaborate on the structure of the Algorithm later, but first provide the main result.

Algorithm 1: The SHRINKING_GRADIENT algorithm.

Data: $T, B > 1, \eta, m$

Result: Weights $\alpha^{(1)}, \dots, \alpha^{(T+1)} \in \mathbb{R}^{(T)}$. Functions $f_t \in L_2(\Omega, \mu)$ defined as

$$f_t = \sum_{i=1}^t \alpha_i^{(t)} \Phi(\mathbf{x}_i);$$

Initialize $\alpha^{(1)} = 0 \in \mathbb{R}^T$;

for $t = 1, \dots, T$ **do**

 Observe \mathbf{x}_t, y_t ;

 Set $E_t = \text{EST_SCALAR_PROD}(\alpha^{(t)}, \mathbf{x}_{1:t-1}, \mathbf{x}_t, m)$;

if $|E_t| < 16B$ **then**

$\alpha^{(t+1)} = \alpha^{(t)}$;

$\alpha_t^{(t+1)} = \eta(E_t - y_t)$;

else

$\alpha^{(t+1)} = \frac{1}{4} \alpha^{(t)}$;

Algorithm 2: EST_SCALAR_PROD

Data: $\alpha, \mathbf{x}_{1:t-1}, \mathbf{x}, m$

Result: Estimated scalar product E

for $k=1 \dots, m$ **do**

 Sample i from the distribution $q(i) = \frac{|\alpha_i|}{\sum |\alpha_i|}$;
 Sample parameter $\bar{\mathbf{w}}$ from μ . Set $E^{(k)} = \text{sgn}(\alpha_i) \psi(\mathbf{x}_i; \bar{\mathbf{w}}) \psi(\mathbf{x}; \bar{\mathbf{w}})$;

Set $E = \frac{\|\alpha\|_1}{m} \sum_{k=1}^m E^{(k)}$

Theorem 1. Run Algorithm 1 with parameters $T, B \geq 1, \eta = \frac{B}{\sqrt{T}}$ and $m = O(B^4 T \log(BT))$. Then:

1. For every sequence of squared losses ℓ_1, \dots, ℓ_T observed by the algorithm we have for f_1, \dots, f_T :

$$\mathbb{E} \left[\sum_{t=1}^T \ell_t(f_t) - \min_{f^* \in \mathcal{H}_\mu^B} \ell_t(f^*) \right] = O(B\sqrt{T})$$

2. The run-time of the algorithm is $\tilde{O}(\rho B^4 T^2)$.³
3. For each $t = 1 \dots T$ and a new test example \mathbf{x} , we can estimate $\langle f_t, \Phi(\mathbf{x}) \rangle$ within accuracy ϵ_0 by running Algorithm 2 with parameters $\alpha^{(t)}, \{\mathbf{x}_i\}_{i=1}^t, \mathbf{x}$ and $m = O(\frac{B^4 T}{\epsilon_0^2} \log 1/\delta)$. The resulting running time at test is then $O(\rho m)$.

We next turn to the statistical setting, where we provide bounds on the expected performance. Following standard online to batch conversion and Theorem 1 we can obtain the following Corollary (e.g., [25]):

Corollary 1 (Statistical Setting). Let $S = \{x_t, y_t\}_{t=1}^T$, be an IID sample drawn from some unknown distribution \mathbb{P} . The following holds for any $\epsilon > 0$. Run Algorithm 1 as in Theorem 1, with $T = O(\frac{B}{\epsilon^2})$. Let $f_S = \frac{1}{T} \sum f_t$. Then the expected loss satisfies:

$$\mathbb{E}_{S \sim \mathbb{P}} [L(f_S)] < \inf_{f^* \in \mathcal{H}_\mu^B} L(f^*) + \epsilon.$$

The runtime of the algorithm, as well as estimation time on a test example are as defined in Theorem 1.

4 Related Work

Learning random features can be traced to the early days of learning [22], and infinite networks have also been introduced more than 20 years ago [29, 17]. More recent works have considered learning neural nets (also multi-layer) with infinite hidden units using

³Ignoring logarithmic factors in B and T .

the kernel trick [8, 13, 15, 16]. These works take a similar approach to ours but focus on computing the kernel for certain feature classes in order to invoke the kernel trick. Our work in contrast avoids using the kernel trick and applies to any feature class that can be randomly generated. All the above works are part of a broader effort of trying to circumvent hardness in deep learning by mimicking deep nets through kernels [21, 4, 2, 3], and developing general duality between neural networks and kernels [11].

From a different perspective the relation between random features and kernels has been noted in [23] where the authors represent translation invariant kernels in terms of random features. This idea has been further studied in [1, 18] for other kernels as well. The focus of these works is mainly to allow scaling down of the feature space and representation of the final output classifier.

The idea is also present in [9], where the authors focus on tractability of large scale kernel methods. More relevant to our work is that one can show that the optimization method in [9] can be invoked whenever the kernel can be estimated using random features. In [9] the objective considered is of the regularized form: $\frac{\gamma}{2}\|f\|^2 + R(f)$, with a corresponding sample complexity of $O(1/(\gamma^2\epsilon^2))$ samples needed to achieve ϵ approximation with respect to the risk of the optimum of the regularized objective.

To relate the above results to ours, we begin by emphasizing that the bound in [9] holds for fixed γ , and refers to optimization of the regularized objective. Our objective is to minimize the risk $R(f)$ which is the expected squared loss, for which we need to choose $\gamma = O(\frac{\epsilon}{B^2})$ in order to attain accuracy ϵ [28]. Plugging this γ into the generalization bound in [9] we obtain that the algorithm in [9] needs $O(\frac{B^4}{\epsilon^4})$ samples to compete with the optimal target function in the B -ball. Our algorithm needs $O(\frac{B}{\epsilon^2})$ examples. We note that their method does extend to a larger class of losses, whereas ours is restricted to the quadratic loss.

In [24], the authors consider embedding the domain into the feature space $\mathbf{x} \rightarrow (\psi(\mathbf{w}_1; \mathbf{x}), \dots, \psi(\mathbf{w}_m; \mathbf{x}))$, where $\mathbf{w}_1, \dots, \mathbf{w}_m$ are IID random variables sampled according to some prior $\mu(\mathbf{w})$. They show that with $O(\frac{B^2 \log 1/\delta}{\epsilon^2})$ random features estimated on $O(\frac{B^2 \log 1/\delta}{\epsilon^2})$ samples they can compete with the class:

$$\mathcal{H}_{\mu_{\max}}^B = \left\{ \mathbf{x} \rightarrow \int \psi(\mathbf{w}; \mathbf{x}) f(\mathbf{w}) d\mu(\mathbf{w}) : |f(\mathbf{w})| \leq B \right\}. \quad (8)$$

Our algorithm relates to the mean square error cost function which does not meet the condition in [24], and is hence formally incomparable. Yet we can invoke our algorithm to compete against a larger class of target functions. Our main result shows that Algorithm 1, using $\tilde{O}(\frac{B^8}{\epsilon^4})$ estimated features and using $O(\frac{B^2}{\epsilon^2})$ samples, will, in expectation, output a predictor that is ϵ close to the best in \mathcal{H}_{μ}^B . Note that $|f(\mathbf{w})| < B$ implies $\mathbb{E}_{\mathbf{w} \sim \mu}(f^2(\mathbf{w})) < B^2$. Hence $\mathcal{H}_{\mu_{\max}}^B \subseteq \mathcal{H}_{\mu}^B$. Note however, that the number of estimated features (as a function of B) is worse in our case.

Our approach to the problem is to consider learning with a noisy estimate of the kernel. A related setting was studied in [7], where the authors considered learning with kernels when the data is corrupted. Noise in the data and noise in the scalar product

estimation are not equivalent when there is non linearity in the kernel space embedding. There is also extensive research on linear regression with actively chosen attributes [6, 14]. The convergence rates and complexity of the algorithms are dimension dependent. It is interesting to see if their method can be extended from finite set of attributes to a continuum set of attributes.

5 Algorithm

We next turn to present Algorithm 1, from which our main result is derived. The algorithm is similar in spirit to Online Gradient Descent (OGD) [31], but with some important modifications that are necessary for our analysis.

We first introduce the problem in the terminology of online convex optimization, as in [31]. At iteration t our algorithm outputs a hypothesis f_t . It then receives as feedback (\mathbf{x}_t, y_t) , and suffers a loss $\ell_t(f_t)$ as in Eq. 5.

The objective of the algorithm is to minimize the regret against a benchmark of B -bounded functions, as in Eq. 6.

A classic approach to the problem is to exploit the OGD algorithm. Its simplest version would be to update $f_{t+1} \rightarrow f_t - \eta \nabla_t$ where η is a step size, and ∇_t is the gradient of the loss wrt f at f_t . In our case, ∇_t is given by:

$$\nabla_t = (y_t - \langle f_t, \Phi(\mathbf{x}_t) \rangle) \Phi(\mathbf{x}_t) \quad (9)$$

Applying this update would also result in a function $f_t = \sum_{i=1}^t \alpha_i \Phi(\mathbf{x}_t)$ as we have in Algorithm 1 (but with different α_i from ours). However, in our setting this update is not applicable since the scalar product $\langle f_t, \Phi(\mathbf{x}_t) \rangle$ is not available.

One alternative is to use a stochastic unbiased estimate of the gradient that we denote by $\bar{\nabla}_t$. This induces an update step $f_{t+1} \rightarrow f_t - \eta \bar{\nabla}_t$. One can show that OGD with such an estimated gradient enjoys the following regret bound for every $\|f^*\| \leq B$ (see for example [25]):

$$\mathbb{E} \left[\sum \ell_t(f_t) - \ell_t(f^*) \right] \leq \frac{B^2}{\eta} + \eta \sum_{i=1}^T \mathbb{E} [\|\nabla_t\|^2] + \eta \sum_{i=1}^T V [\bar{\nabla}_t]. \quad (10)$$

Where $V [\bar{\nabla}_t] = \mathbb{E} [\|\bar{\nabla}_t - \nabla_t\|^2]$. We can bound the first two terms using standard techniques applicable for the squared loss (as appeared for example in [30] or [27]). The third term depends on our choice of gradient estimate. There are multiple possible choices for such an estimate, and we use a version which facilitates our analysis.

Assume that at iteration t , our function f_t is given by $f_t = \sum_{i=1}^t \alpha_i^{(t)} \Phi(\mathbf{x}_t)$. We now want to use sampling to obtain an unbiased estimate of $\langle f_t, \Phi(\mathbf{x}_t) \rangle$. We will do this via a two step sampling procedure, as described in Algorithm 2. First, sample an index $i \in [1, \dots, t]$ by sampling according to the distribution $q(i) \propto |\alpha_i^{(t)}|$. Next, for the chosen i , sample $\bar{\mathbf{w}}$ according to μ , and use $\psi(\mathbf{x}; \bar{\mathbf{w}})\psi(\mathbf{x}_i; \bar{\mathbf{w}})$ to construct an estimate of

$\langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_t) \rangle$. The resulting unbiased estimate of $\langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_t) \rangle$ is denoted by E_t and given by:

$$E_t = \frac{\|\alpha^{(t)}\|_1}{m} \sum_{i=1}^m \text{sgn}(\alpha_i^{(t)}) \psi(\mathbf{x}_i; \bar{\mathbf{w}}) \psi(\mathbf{x}_t; \bar{\mathbf{w}}) \quad (11)$$

The corresponding unbiased gradient estimate is:

$$\bar{\nabla}_t = (E_t - y_t) \mathbf{x}_t \quad (12)$$

The variance of $\bar{\nabla}$ affects the convergence rate and depends on both $\|\alpha\|_1$ and m —the number of estimations. We wish to maintain $m = O(T)$ estimations per round, while achieving $O(\sqrt{T})$ regret.

To effectively regularize $\|\alpha\|_1$, we modify the OGD algorithm so that whenever E_t is larger than $16B$, we do not perform the usual update. Instead, we perform a shrinking step that divides $\alpha^{(t)}$ (and hence f_t) by 4. Treating B as constant, this guarantees that $\|\alpha\|_1 = O(\eta T)$, and in turn $\text{Var}(\bar{\nabla}_t) = O(\frac{\eta^2 T^2}{m})$. Setting $\eta = O(1/\sqrt{T})$, we have that $m = O(T)$ estimations are enough.

The rationale behind the shrinkage is that whenever E_t is large, it indicates that f_t is “far away” from the B -ball, and a shrinkage step, similar to projection, leads f_t closer to the optimal element in the B -ball hence improves f_t . However, due to stochasticity, the shrinkage step does add a further term to the regret bound that we would need to take care of.

5.1 Analysis

In what follows we analyze the regret for Algorithm 1. We begin by modifying the regret bound for OGD in Eq. 10 to accommodate for steps that differ from the standard gradient update, such as shrinkage.

Lemma 1. *Let ℓ_1, \dots, ℓ_T be an arbitrary sequence of convex loss functions, and let f_1, \dots, f_T be random vectors, produced by an online algorithm. Assume $\|f_i\| \leq B_T$ for all $i \leq T$. For each t let $\bar{\nabla}_t$ be an unbiased estimator of $\nabla \ell_t(f_t)$. Denote $\hat{f}_t = f_{t-1} - \eta \bar{\nabla}_{t-1}$ and let*

$$P_t(f^*) = \mathbb{P} \left[\|f_t - f^*\| > \|\hat{f}_t - f^*\| \right]. \quad (13)$$

For every $\|f^*\| \leq B$ it holds that:

$$\mathbb{E} \left[\sum_{t=1}^T \ell_t(f_t) - \ell_t(f^*) \right] \leq \frac{B^2}{\eta} + \eta \sum_{t=1}^T \mathbb{E} [\|\nabla_t\|^2] + \eta \sum_{t=1}^T V[\bar{\nabla}_t] + \sum_{t=1}^T \frac{(B_T + B)^2}{\eta} \mathbb{E} [P_{t+1}(f^*)] \quad (14)$$

As discussed earlier, the first two terms in the RHS are the standard bound for OGD from Eq. 10. Note that in OGD we always choose $f_t = \hat{f}_t$ therefore $P_t(f^*) = 0$ and the last term disappears.

The third term will be bounded by controlling $\|\alpha\|_1$. The last term $P_{t+1}(f^*)$ is a penalty that results from updates that stir f_t away from the standard update step \hat{f}_t . This will indeed happen for the shrinkage step. The next lemma bounds this term.

Lemma 2. *Run Algorithm 1 with parameters T , $B \geq 1$ and $\eta < 1/8$. Let $\bar{\nabla}_t$ be the unbiased estimator of $\nabla \ell_t(f_t)$ of the form*

$$\bar{\nabla}_t = (E_t - y_t)\Phi(\mathbf{x}_t).$$

Denote $\hat{f}_t = f_t - \eta \bar{\nabla}_t$ and define $P_t(f^)$ as in Eq. 13.*

$$P_t(f^*) \leq 2 \exp \left(-\frac{m}{(3\eta t)^2} \right)$$

Considering Lemma 2 and Lemma 1, we obtain that the last term in Eq. 14 can be bounded by ηT if we consider $m = O \left((\eta T)^2 \log \left(\frac{B_T}{\eta} \right)^2 \right)$ estimations (while again for simplicity, treating B as constant). One can show that $B_T = O(\eta T)$, hence neglecting logarithmic factors and by choice of $\eta = O(\frac{1}{\sqrt{T}})$, we obtain that $m = \tilde{O}(T)$ estimations per round suffice to bound the fourth term. A full proof of Theorem 1 and the above two Lemmas is given in the Appendix.

6 Discussion

We presented a new online algorithm that employs kernels implicitly but avoids the kernel trick assumption. Namely, the algorithm can be invoked even when one has access to only estimations of the scalar product. The problem was motivated by kernels resulting from neural nets, but it can of course be applied to any scalar product of the form we described. To summarize, we can cast our result into a larger model that might be of independent interest. Consider a setting where a learner can observe an unbiased estimate of a coordinate in a kernel matrix, or alternatively the scalar product between any two observations. Our results imply that in this setting the above rates are applicable, and at least for the square loss, having no access to the true values in the kernel matrix is not necessarily prohibitive during training.

The results show that with sample size T we can achieve error of $O(\frac{B}{\sqrt{T}})$. As demonstrated by [26] these rates are optimal, even when the scalar product is computable. To achieve this rate our algorithm needs to perform $\tilde{O}(B^4 T^2)$ scalar product estimations. When the scalar product can be computed, existing kernelized algorithms need to observe a fixed proportion of the kernel matrix, hence they observe order of $\Omega(T^2)$ scalar products. In [5] it was shown that when the scalar product can be computed exactly, one would need access to at least $\Omega(T)$ entries to the kernel matrix. It is still an open problem whether one has to access $\Omega(T^2)$ entries when the kernel can be computed exactly. However, as we show here, for fixed B even if the kernel can only be estimated $\tilde{O}(T^2)$ estimations are enough. It would be interesting to further investigate and improve the performance of our algorithm in terms of the norm bound B .

Acknowledgement

The authors would like to thank Tomer Koren for helpful comments and suggestions. Roi Livni is a recipient of the Google Europe Fellowship in Learning Theory, and this research is supported in part by this Google Fellowship, and a Google Research Award.

References

- [1] F. Bach. On the equivalence between kernel quadrature rules and random feature expansions. 2015.
- [2] L. Bo, K. Lai, X. Ren, and D. Fox. Object recognition with hierarchical kernel descriptors. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 1729–1736. IEEE, 2011.
- [3] L. Bo, X. Ren, and D. Fox. Kernel descriptors for visual recognition. In *Advances in neural information processing systems*, pages 244–252, 2010.
- [4] J. Bouvrie, L. Rosasco, and T. Poggio. On invariance in hierarchical models. In *Advances in Neural Information Processing Systems*, pages 162–170, 2009.
- [5] N. Cesa-Bianchi, Y. Mansour, and O. Shamir. On the complexity of learning with kernels. In *Proceedings of The 28th Conference on Learning Theory*, pages 297–325, 2015.
- [6] N. Cesa-Bianchi, S. Shalev-Shwartz, and O. Shamir. Efficient learning with partially observed attributes. *The Journal of Machine Learning Research*, 12:2857–2878, 2011.
- [7] N. Cesa-Bianchi, S. Shalev-Shwartz, and O. Shamir. Online learning of noisy data. *Information Theory, IEEE Transactions on*, 57(12):7907–7931, 2011.
- [8] Y. Cho and L. K. Saul. Kernel methods for deep learning. In *Advances in neural information processing systems*, pages 342–350, 2009.
- [9] B. Dai, B. Xie, N. He, Y. Liang, A. Raj, M. Balcan, and L. Song. Scalable kernel methods via doubly stochastic gradients. *CoRR*, abs/1407.5599, 2014.
- [10] B. Dai, B. Xie, N. He, Y. Liang, A. Raj, M.-F. F. Balcan, and L. Song. Scalable kernel methods via doubly stochastic gradients. In *Advances in Neural Information Processing Systems*, pages 3041–3049, 2014.
- [11] A. Daniely, R. Frostig, and Y. Singer. Toward deeper understanding of neural networks: The power of initialization and a dual view on expressivity. *arXiv preprint arXiv:1602.05897*, 2016.
- [12] A. Daniely, N. Linial, and S. Shalev-Shwartz. From average case complexity to improper learning complexity. In *Proceedings of the 46th Annual ACM Symposium on Theory of Computing*, pages 441–448. ACM, 2014.
- [13] L. Deng, G. Tur, X. He, and D. Hakkani-Tur. Use of kernel deep convex networks and end-to-end learning for spoken language understanding. In *Spoken Language Technology Workshop (SLT), 2012 IEEE*, pages 210–215. IEEE, 2012.
- [14] E. Hazan and T. Koren. Linear regression with limited observation. In *Proceedings of the 29th International Conference on Machine Learning (ICML-12)*, pages 807–814, 2012.

- [15] T. Hazan and T. Jaakkola. Steps toward deep kernel methods from infinite neural networks. *arXiv preprint arXiv:1508.05133*, 2015.
- [16] U. Heinemann, R. Livni, E. Eban, G. Elidan, and A. Globerson. Improper deep kernels. In *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, pages 1159–1167, 2016.
- [17] K. Hornik. Some new results on neural network approximation. *Neural Networks*, 6(8):1069–1072, 1993.
- [18] P. Kar and H. Karnick. Random feature maps for dot product kernels. In *International Conference on Artificial Intelligence and Statistics*, pages 583–591, 2012.
- [19] A. R. Klivans and A. A. Sherstov. Cryptographic hardness for learning intersections of half-spaces. In *Foundations of Computer Science, 2006. FOCS'06. 47th Annual IEEE Symposium on*, pages 553–562. IEEE, 2006.
- [20] R. Livni, S. Shalev-Shwartz, and O. Shamir. On the computational efficiency of training neural networks. In *Advances in Neural Information Processing Systems*, pages 855–863, 2014.
- [21] J. Mairal, P. Koniusz, Z. Harchaoui, and C. Schmid. Convolutional kernel networks. In *Advances in Neural Information Processing Systems*, pages 2627–2635, 2014.
- [22] M. Minsky and S. Papert. Perceptrons: an introduction to computational geometry (expanded edition), 1988.
- [23] A. Rahimi and B. Recht. Random features for large-scale kernel machines. In *Advances in neural information processing systems*, pages 1177–1184, 2007.
- [24] A. Rahimi and B. Recht. Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning. In *Advances in neural information processing systems*, pages 1313–1320, 2009.
- [25] S. Shalev-Shwartz. Online learning and online convex optimization. *Foundations and Trends in Machine Learning*, 4(2):107–194, 2011.
- [26] O. Shamir. The sample complexity of learning linear predictors with the squared loss. *Journal of Machine Learning Research*, 16(Dec):3475–3486, 2015.
- [27] N. Srebro, K. Sridharan, and A. Tewari. Smoothness, low noise and fast rates. In *Advances in neural information processing systems*, pages 2199–2207, 2010.
- [28] K. Sridharan, S. Shalev-Shwartz, and N. Srebro. Fast rates for regularized objectives. In *Advances in Neural Information Processing Systems*, pages 1545–1552, 2009.
- [29] C. K. Williams. Computing with infinite networks. *Advances in neural information processing systems*, pages 295–301, 1997.
- [30] T. Zhang. Solving large scale linear prediction problems using stochastic gradient descent algorithms. In *Proceedings of the twenty-first international conference on Machine learning*, page 116. ACM, 2004.
- [31] M. Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Machine Learning, Proceedings of the Twentieth International Conference (ICML)*, pages 928–936, 2003.

A Estimation Procedure – Concentration Bounds

In this section we provide concentration bounds for the estimation procedure in Algorithm 2.

Lemma 3. *Run Algorithm 2 with α and, $\{\mathbf{x}_i\}_{i=1}^T$, and m . Let $f = \sum \alpha_i \Phi(\mathbf{x}_i)$. Assume that $|\psi(\mathbf{x}; \mathbf{w})| < 1$ for all \mathbf{w} and \mathbf{x} . Let E be the output of Algorithm 2. Then E is an unbiased estimator for $\langle f, \Phi(\mathbf{x}) \rangle$ and:*

$$\mathbb{P}[|E - \langle f, \Phi(\mathbf{x}) \rangle| > \epsilon] \leq \exp\left(-\frac{m\epsilon^2}{\|\alpha\|_1^2}\right) \quad (15)$$

Proof. Consider the random variables $\|\alpha\|_1 E^{(k)}$ (where $E^{(k)}$ is as defined in Algorithm 2) and note that they are IID. One can show that $\mathbb{E}[\|\alpha\|_1 E^{(k)}] = \sum \alpha_i \mathbb{E}[\psi(\Phi(\mathbf{x}_i); \bar{\mathbf{w}}) \psi(\Phi(\mathbf{x}); \bar{\mathbf{w}})] = \langle f, \Phi(\mathbf{x}) \rangle$. By the bound on $\psi(\mathbf{x}; \mathbf{w})$ we have that $|\|\alpha\|_1 E^{(k)}| < \|\alpha\|_1$ with probability 1. Since $E = \frac{1}{m} \sum E^{(k)}$ the result follows directly from Hoeffding's inequality. \square

Next, we relate these guarantees to the output of Algorithm 1:

Lemma 4. *The $\alpha^{(t)}$ obtained in Algorithm 1 satisfies:*

$$\|\alpha^{(t)}\|_1 \leq (16B + 1)\eta t.$$

As a corollary of this and Lemma 3 we have that the function f_t satisfies:

$$\mathbb{P}[|E_t - \langle f_t, \Phi(\mathbf{x}_t) \rangle| > \epsilon] \leq \exp\left(-\frac{\epsilon^2 m}{((16B + 1)\eta t)^2}\right) \quad (16)$$

Proof. We prove the statement by induction. We separate into two cases, depending on whether the shrinkage step was performed or not.

If $|E_t| \geq 16B$ the algorithm sets $\alpha^{(t+1)} = \frac{1}{4}\alpha^{(t)}$, and:

$$\|\alpha^{(t+1)}\|_1 = \frac{1}{4}\|\alpha^{(t)}\|_1 \leq (16B + 1)\eta(t + 1)$$

If $|E_t| < 16B$ the gradient update is performed. Since $|y_t| \leq 1$ we have that $|E_t - y_t| < 16B + 1$ and:

$$\|\alpha^{(t+1)}\|_1 \leq \|\alpha^{(t)}\|_1 + \eta|E_t - y_t| \leq (16B + 1)\eta(t + 1).$$

\square

B Proofs of Lemmas

B.1 Proof of Lemma 1

First, by convexity we have that

$$2(\ell_t(f_t) - \ell_t(f^*)) \leq 2\langle \nabla_t, f_t - f^* \rangle. \quad (17)$$

Next we upper bound $\langle \nabla_t, f_t - f^* \rangle$. Denote by \mathcal{E} the event $\|f_{t+1} - f^*\| > \|\hat{f}_{t+1} - f^*\|$. Note that

$$\begin{aligned} \mathbb{E} [\|f_{t+1} - f^*\|^2] &\leq \mathbb{E} [\|\hat{f}_{t+1} - f^*\|^2] + \mathbb{E} [\|f_{t+1} - f^*\|^2 | \mathcal{E}] \cdot P_{t+1}(f^*) \leq \\ &\mathbb{E} [\|\hat{f}_{t+1} - f^*\|^2] + (B + B_T)^2 P_{t+1}(f^*) \end{aligned}$$

Plugging in $\hat{f}_{t+1} = f_t - \eta \bar{\nabla}_t$ we get

$$\mathbb{E} [\|f_{t+1} - f^*\|^2] \leq \mathbb{E} [\|f_t - f^*\|^2] + \eta^2 \mathbb{E} [\|\bar{\nabla}_t\|^2] - 2\eta \mathbb{E} [\langle \bar{\nabla}_t, f_t - f^* \rangle] + (B + B_T)^2 P_{t+1}(f^*)$$

Dividing by η we have that:

$$2\mathbb{E} [\langle \bar{\nabla}_t, f_t - f^* \rangle] \leq \frac{\mathbb{E} [\|f_t - f^*\|^2] - \mathbb{E} [\|f_{t+1} - f^*\|^2]}{\eta} + \frac{1}{\eta} \mathbb{E} [(B + B_T)^2 P_{t+1}(f^*)] + \eta \mathbb{E} [\|\bar{\nabla}_t\|^2] \quad (18)$$

Taking 17 and 18 and summing we have:

$$\begin{aligned} 2\mathbb{E} \left[\sum_{t=1}^T \ell_t(f_t) - \ell_t(f^*) \right] &\leq 2\mathbb{E} [\nabla_t^\top (f_t - f^*)] = 2\mathbb{E} [\bar{\nabla}_t^\top (f_t - f^*)] \leq \\ &\frac{\|f^*\|^2}{\eta} + \frac{1}{\eta} \sum_{t=1}^T \mathbb{E} [(B + B_T)^2 P_{t+1}(f^*)] + \eta \sum_{t=1}^T \mathbb{E} [\|\bar{\nabla}_t\|^2]. \end{aligned}$$

Finally note that $\mathbb{E} [\|\bar{\nabla}_t\|^2] = \mathbb{E} [\|\nabla_t\|^2] + V [\bar{\nabla}_t]$ to obtain the result.

B.2 Proof for Lemma 2

To prove the bound in the lemma, we first bound the event $P_t(f^*)$ w.r.t to two possible events:

Lemma 5. *Consider the setting as in Lemma 2. Run Algorithm 1 and for each t consider the following two events:*

- \mathcal{E}_1^t : $|E_t| > 16B$ and $|E_t| > \frac{1}{4\eta} \|f_t\|$.

- \mathcal{E}_2^t : $|E_t| > 16B$ and $\|f_t\| < 8B$.

For every $\|f^*\| < B$ we have that $P_t(f^*) < \mathbb{P} [\mathcal{E}_1^t \cup \mathcal{E}_2^t]$.

Proof. Denote the event $|E_t| > 16B$ by \mathcal{E}_0^t . Note that if \mathcal{E}_0^t does not happen, then $f_t = \hat{f}_t$. Hence trivially

$$P_t(f^*) = \mathbb{P} [\|f_t - f^*\| > \|\hat{f}_t - f^*\| \wedge \mathcal{E}_0^t]$$

We will assume that:

1. $|E_t| > 16B$.
2. $|E_t| < \frac{1}{4\eta}\|f_t\|$.
3. $\|f_t\| > 8B$

and show $\|f_{t+1} - f^*\| \leq \|\hat{f}_{t+1} - f^*\|$.

In other words, we will show that if \mathcal{E}_0^t happens and $\|f_{t+1} - f^*\| > \|\hat{f}_{t+1} - f^*\|$, then either \mathcal{E}_2^t or \mathcal{E}_1^t happened. This will conclude the proof.

Fix t , note that since $|\psi(\mathbf{x}; \mathbf{w})| < 1$ we have that $\|\Phi(\mathbf{x})\| < 1$. We then have:

$$\|\hat{f}_{t+1}\| = \|f_t - \eta(E_t - y)\Phi(\mathbf{x}_t)\| > \|f_t\| - \|\eta(E_t - y)\Phi(\mathbf{x}_t)\| \geq \|f_t\| - \eta|E_t| - \eta \geq \frac{3}{4}\|f_t\| - \eta$$

Where the last inequality is due to assumption 2. We therefore have the following bound for every $\|f^*\| < B$:

$$\|\hat{f}_{t+1} - f^*\| \geq \frac{3}{4}\|f_t\| - \eta - B$$

On the other hand, if $f_{t+1} \neq \hat{f}_{t+1}$ then by construction of the algorithm $f_{t+1} = \frac{1}{4}f_t$:

$$\|f_{t+1} - f^*\| \leq \|f_{t+1}\| + \|f^*\| \leq \frac{\|f_t\|}{4} + B.$$

Next note that $\eta < 2B$ and assumption 3 states $\|f_t\| > 8B$. Therefore:

$$\frac{1}{2}\|f_t\| > 4B > \eta + 2B$$

and we obtained the desired result:

$$\|\hat{f}_{t+1} - f^*\| \geq \frac{3}{4}\|f_t\| - \eta - B = \frac{1}{4}\|f_t\| + \left(\frac{1}{2}\|f_t\| - \eta - 2B\right) + B \geq \frac{1}{4}\|f_t\| + B \geq \|f_{t+1} - f^*\|$$

□

Next we upper bound $\mathbb{P}[\mathcal{E}_1^t \cup \mathcal{E}_2^t]$. In what follows the superscript t is dropped.

A bound for $\mathbb{P}[\mathcal{E}_1 \cap \mathcal{E}_2^c]$: Assume that

$$|E_t - \langle f_t, \Phi(\mathbf{x}_t) \rangle| < \left(\frac{1}{4\eta} - 1\right)8B.$$

We assume T is sufficiently large and $\eta < \frac{1}{8}$. We have $\frac{1}{4\eta} - 1 > 1$.

Since we assume \mathcal{E}_2 did not happen we must have $\|f_t\| > 8B$ and

$$|E_t - \langle f_t, \Phi(\mathbf{x}_t) \rangle| < \left(\frac{1}{4\eta} - 1\right)\|f\|.$$

We continue:

$$E_t - \|f\| < |E_t - \langle f_t, \Phi(\mathbf{x}_t) \rangle| < \left(\frac{1}{4\eta} - 1\right)\|f\|.$$

Which leads to

$$E_t < \frac{1}{4\eta} \|f\|.$$

And we get that \mathcal{E}_1 did not happen. We conclude that if \mathcal{E}_1 and not \mathcal{E}_2 then:

$$|E_t - \langle f_t, \Phi(\mathbf{x}_t) \rangle| \geq (\frac{1}{4\eta} - 1)8B.$$

Since $\frac{1}{4\eta} - 1 > 1$ we have that:

$$|E_t - \langle f_t, \Phi(\mathbf{x}_t) \rangle| \geq 8B.$$

We conclude that:

$$\mathbb{P}[\mathcal{E}_1 \cap \mathcal{E}_2^c] \leq \mathbb{P}[|E_t - \langle f_t, \Phi(\mathbf{x}_t) \rangle| \geq 8B]. \quad (19)$$

A bound for $\mathbb{P}[\mathcal{E}_2]$: If $|E_t| > 16B$ and $\|f_t\| < 8B$ then by normalization of $\Phi(\mathbf{x}_t)$ we have that $\langle f_t, \Phi(\mathbf{x}_t) \rangle < 8B$ and trivially we have that

$$|E_t - \langle f_t, \Phi(\mathbf{x}_t) \rangle| \geq 8B.$$

And again we have that:

$$\mathbb{P}[\mathcal{E}_2] \leq \mathbb{P}[|E_t - \langle f_t, \Phi(\mathbf{x}_t) \rangle| \geq 8B]. \quad (20)$$

Taking Eq. 19 and Eq. 20 we have that

$$\mathbb{P}[\mathcal{E}_2 \cup \mathcal{E}_1] \leq 2\mathbb{P}[|E_t - \langle f_t, \Phi(\mathbf{x}_t) \rangle| \geq 8B]. \quad (21)$$

By Lemma 4 we have that:

$$\begin{aligned} &P(|E_t - \langle f_t, \Phi(\mathbf{x}_t) \rangle| > 8B) < \\ &\exp\left(-\frac{m(8B)^2}{((16B+1)\eta t)^2}\right) < \exp\left(-\frac{m}{(3\eta t)^2}\right) \end{aligned}$$

Taking the above upper bounds together with Lemma 5 we can prove Lemma 2.

C Proof of Main Result

C.1 Some more technical Lemmas

We begin by deriving Corollary 2 that bounds the first two terms in the regret bound. As discussed, this section follows standard techniques. We begin with an upper bound on $\mathbb{E}(\|\bar{\nabla}_t\|^2)$.

Lemma 6. *Consider the setting as in Lemma 2. Then*

$$V[\bar{\nabla}_t] \leq \frac{((16B+1)\eta t)^2}{m},$$

and,

$$\mathbb{E}[\|\nabla_t\|^2] \leq 2\mathbb{E}[\ell_t(f_t)].$$

Proof. Begin by noting that since $\|\Phi(\mathbf{x})\| < 1$, it follows from the definitions of $\nabla, \bar{\nabla}$ that:

$$V[\bar{\nabla}_t] = \mathbb{E}[\|\bar{\nabla}_t - \nabla_t\|^2] \leq \mathbb{E}[(E_t - \langle f_t, \Phi(\mathbf{x}_t) \rangle)^2] = V[E_t]$$

By construction (see Algorithm 2) we have that:

$$V[E_t] = \frac{1}{m} V[\|\alpha^{(t)}\|_1^2 \psi(\mathbf{x}_i; \mathbf{w}) \psi(\mathbf{x}_t; \mathbf{w})]$$

where the index i is sampled as in Algorithm 2, and $\psi(\mathbf{x}_i; \mathbf{w}) \psi(\mathbf{x}_t; \mathbf{w})$ is bounded by 1. By Lemma 4 we have that

$$V[E_t] \leq \frac{((16B+1)\eta t)^2}{m}.$$

This provides the required bound on $V[\bar{\nabla}_t]$. Additionally, we have that

$$\|\nabla_t\|^2 = (\langle f_t, \Phi(\mathbf{x}_t) \rangle - y_t)^2 \|\Phi(\mathbf{x}_t)\|^2 \leq 2\ell_t(f_t)$$

and the result follows by taking expectation. \square

We thus have the following corollary that bounds the first three terms in Eq. 14:

Corollary 2. *With the notations and setting of Lemma 2 we have:*

$$\begin{aligned} \frac{B^2}{\eta} + \eta \sum_{t=1}^T \mathbb{E}[\|\nabla_t\|^2] + \eta \sum_{t=1}^T V[\bar{\nabla}_t] &\leq 2\eta \mathbb{E} \left[\sum_{t=1}^T \ell_t(f_t) - \ell_t(f^*) \right] + \frac{B^2}{\eta} \\ &\quad + 2\eta \sum_{t=1}^T \ell_t(f^*) + \eta \sum_{t=1}^T \frac{((16B+1)\eta T)^2}{m} \end{aligned} \quad (22)$$

C.2 Proof of Theorem 1

Recall that Lemma 1 and Corollary 2 assume an upper bound B_T on $\|f_t\|$. We begin by noting that B_T can be bounded as follows, using Lemma 4:

$$B_T = \max_t \|f_t\| \leq \max_t \|\alpha^{(t)}\|_1 \leq (16B+1)\eta T. \quad (23)$$

Plugging Corollary 2 into Eq. 14 we obtain:

$$(1-2\eta) \mathbb{E} \left[\sum_{t=1}^T \ell_t(f_t) - \ell_t(f^*) \right] \leq \frac{B^2}{\eta} + 2\eta \sum_{t=1}^T \ell_t(f^*) + \eta \sum_{t=1}^T \frac{((16B+1)\eta T)^2}{m} + \frac{(B_T + B)^2}{\eta} \sum_{t=1}^T P_t(f^*) \quad (24)$$

To bound the second term we note that:

$$\min_{\|f^*\| < B} \sum_{t=1}^T \ell_t(f^*) \leq \sum_{t=1}^T \ell_t(0) \leq T. \quad (25)$$

We next set η and m as in the statement of the theorem. Namely: $\eta = \frac{B}{2\sqrt{T}}$, and $m = ((16B+1)B)^2 T \log \gamma$, where $\gamma = \max\left(\frac{((16B+1)\eta T+B)^2}{\eta^2}, e\right)$.

Our choice of m implies that $m > ((16B+1)\eta T)^2$, and hence the third term in Eq. 24 is bounded as follows:

$$\eta \sum_{t=1}^T \frac{((16B+1)\eta T)^2}{m} \leq \eta T \quad (26)$$

Next we have that $m > (3\eta t)^2 \log \gamma$ for every t , and by the bound on B_T we have that $\gamma > \frac{(B+B_T)^2}{\eta^2}$. Taken together with Lemma 2 we have that:

$$\frac{(B_T+B)^2}{\eta} \sum_{t=1}^T P_t(f^*) \leq \eta T. \quad (27)$$

Taking Eq. 25, Eq. 27 and Eq. 26, plugging them into Eq. 24 we have that:

$$(1-2\eta)\mathbb{E} \left[\sum_{t=1}^T \ell_t(f_t) - \ell_t(f^*) \right] \leq \frac{B^2}{\eta} + 2\eta T + \eta T + \eta T$$

Finally by choice of η , and dividing both sides by $(1-2\eta)$ we obtain the desired result.

It remains to show that we can estimate each f_t in the desired complexity (the result for the averaged f is the same). Each f_t has the form $f_t = \sum_{i=1}^T \alpha_i^{(t)} \mathbf{x}_i$, By Lemma 3 and Lemma 4, running Algorithm 2 m iterations will lead to a random variable E such that:

$$\mathbb{P} [|E - \langle f_t, \mathbf{x} \rangle|] \leq \exp \left(-\frac{\epsilon^2 m}{((16B+1)B\sqrt{T})^2} \right).$$

We obtain that order of $O(\frac{B^4 T}{\epsilon^2} \log 1/\delta)$ estimations are enough.