
Tracking the Best Regressor*

Mark Herbster and Manfred K. Warmuth

Department of Computer Science
University of California at Santa Cruz
Applied Sciences Building
Santa Cruz, CA 95064 (USA)
mark@cs.ucsc.edu, manfred@cs.ucsc.edu

Abstract

In most of the on-line learning research the total on-line loss of the algorithm is compared to the total loss of the best off-line predictor u from a comparison class of predictors. We call such bounds static bounds. The interesting feature of these bounds is that they hold for an arbitrary sequence of examples. Recently some work has been done where the comparison vector u_t at each trial t is allowed to change with time, and the total on-line loss of the algorithm is compared to the sum of the losses of u_t at each trial plus the total “cost” for shifting to successive comparison vectors. This is to model situations in which the examples change over time and different predictors from the comparison class are best for different segments of the sequence of examples. We call such bounds shifting bounds. Shifting bounds still hold for arbitrary sequences of examples and also for arbitrary partitions. The algorithm does not know the off-line partition and the sequence of predictors that its performance is compared against.

Naturally shifting bounds are much harder to prove. The only known bounds are for the case when the comparison class consists of a finite sets of experts or boolean disjunctions. In this paper we develop the methodology for lifting known static bounds to the shifting case. In particular we obtain bounds when the comparison class consists of linear neurons (linear combinations of experts). Our essential technique consists of the following. At the end of each trial we *project* the hypothesis of the static algorithm into a suitably chosen convex region. This keeps the hypothesis of the algorithm well-behaved and the static bounds can be converted to shifting bounds so that the cost for shifting remains reasonable.

*The authors were supported by the NSF grant CCR-9700201.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

COLT 98 Madison WI USA

Copyright ACM 1998 1-58113-057-0/98/ 7...\$5.00

1 INTRODUCTION

Consider the following by now standard on-line learning model which is a generalization of a model introduced by Littlestone [Lit89, Lit88, LW94]. Learning proceeds in trials $t = 1, 2, \dots, \ell$. The algorithm maintains a parameter vector (hypothesis), denoted by $w_t \in \mathbb{R}^n$. In each trial the algorithm receives an *instance* x_t . It then produces some action or a prediction for x_t based on w_t . Finally the algorithm receives an *outcome* y_t and incurs a loss describing how well the hypothesis w_t “performed” on the *example* (x_t, y_t) . This loss is some non-negative function $L(w_t, (x_t, y_t))$.

For example, the algorithm might predict with $\sigma(w_t \cdot x_t)$, where σ is the logistic function, and the loss might be the square loss, i.e. $\frac{1}{2}(\sigma(w_t \cdot x_t) - y_t)^2$. Thus in this example w represents a hypothesis $h_w(x) = \sigma(w \cdot x)$. By choosing a different function σ , the hypotheses class of the algorithm changes. We abbreviate the loss in trial t by $L_t(w_t)$ and the total loss of the algorithm A on a sequence $S = \langle (x_1, y_1), \dots, (x_\ell, y_\ell) \rangle$ of examples by $L(A, S) = \sum_{t=1}^{\ell} L_t(w_t)$.

In the methodology of worst-case loss bounds the total loss of the algorithm is expressed as a function of the total loss of any member in a comparison class of predictors, which is usually a subset of the hypotheses class. Such a predictor u also incurs a loss $L_t(u)$ in each trial and the total loss of u on the entire sequence is abbreviated as $L(u, S) = \sum_{t=1}^{\ell} L_t(u)$. In the simplest case the bounds have the following form

$$L(A, S) \leq cL(u, S) + c \text{size}(u). \quad (1)$$

Here c is a small constant, u is any vector in the comparison class, and $\text{size}(u)$ is a measurement of the size of u . We call such bounds *static* bounds, because the comparison vector u is any member from the comparison class but it does not change with time. Surprisingly, such bounds are achievable even when there are no probabilistic assumptions made on the sequence of examples [Lit88, Vov90, HKW97, CBLW96, KW94, Byl97, HKW95]. In this paper we allow the comparison vector u to shift with time. For a sequence S of examples of length ℓ and a schedule of predictors $\langle u_1, \dots, u_\ell \rangle$ from the comparison class, we seek an upper bound of the form

$$L(A, S) \leq cL(\langle u_1, \dots, u_\ell \rangle, S) + c \text{size}(\langle u_1, \dots, u_\ell \rangle). \quad (2)$$

Here $L(\langle \mathbf{u}_1, \dots, \mathbf{u}_\ell \rangle, S)$ is the loss of the schedule of predictors on S , i.e. $L(\langle \mathbf{u}_1, \dots, \mathbf{u}_\ell \rangle, S) = \sum_{t=1}^{\ell} L_t(\mathbf{u}_t)$, and, intuitively, $\text{size}(\langle \mathbf{u}_1, \dots, \mathbf{u}_\ell \rangle)$ measures the amount of shifting that occurs in the schedule. We call such bounds *shifting bounds*. In this paper our bounds use the following measure of $\text{size}(\langle \mathbf{u}_1, \dots, \mathbf{u}_\ell \rangle)$:

$$\|\langle \mathbf{u}_1, \dots, \mathbf{u}_\ell \rangle\|_p = \sum_{t=1}^{\ell-1} \|\mathbf{u}_t - \mathbf{u}_{t+1}\|_p + \|\mathbf{u}_\ell\|_p,$$

where $\|\cdot\|_p$ is a p-norm. This is in accord with previous work on worst case loss bounds, where the static bound also grows with a p-norm of the comparison vector. The new shifting bounds presented in this paper build on previous work of the authors [HW98] where the loss of the algorithm was compared against the loss of the best shifting expert (see also [Vov97]) or the best shifting disjunction [AW98]. The work on shifting experts has been applied to predicting disk idle times [HLS96] as well as load balancing problems [BB97]. In this paper we discuss general methods of lifting static bounds to the more complicated shifting case.

We focus on a class of algorithms that is characterized by a *link function* f . Such functions are arbitrary mappings for which the domain is some open set in \mathbb{R}^n , the range is a subset of \mathbb{R}^n , and furthermore the Jacobian of f is strictly positive definite on the whole domain. The update for a given link function f is as follows:

$$\mathbf{w}_{t+1} = f^{-1}(f(\mathbf{w}_t) - \eta \nabla L_t(\mathbf{w}_t)). \quad (3)$$

Here, f^{-1} denotes the inverse of f and η is a non-negative learning rate. The general form of these updates was introduced in [GLS97] for linear threshold functions and in [JW98, KW97b] for regression problems. We refer to these updates as *generalized gradient descent algorithms*. For the sake of simplicity we focus on only two link functions in this paper, the identity function and the logarithm function. For the identity link $\text{id}(\mathbf{w}) = \mathbf{w}$ the domain is all of \mathbb{R}^n and the above update becomes the standard Gradient Descent Algorithm (GD). For the logarithm link function the domain is $(0, \infty)^n$ and $\ln(\mathbf{w})$ is defined as the component-wise application of the \ln function, i.e. $\ln(\mathbf{w}) = (\ln w_1, \dots, \ln w_n)$. The generalized gradient descent algorithm for the \ln link is the Un-normalized Exponentiated Gradient (EGU) Algorithm [KW97a] (see Figure 1).

The key tool for analyzing the update associated with a link function f is a divergence function introduced by Bregman [Bre67, CL81, Csi91]. These divergences were previously used in convex programming in connection with generalized maximum entropy methods. In the context of on-line learning various versions of these divergence functions were first used in [KW97b, GLS97, JW98]). At this point we introduce the notion of divergence without listing a number of technical conditions specified in the Appendix. Given a strictly convex function $F : E \rightarrow \mathbb{R}$, where $E \subset \mathbb{R}^n$ is a closed convex set and $\text{ri } E$ denotes the relative interior of E , then the Bregman divergence $D_F : E \times \text{ri } E \rightarrow [0, \infty)$ is defined as

$$D_F(\mathbf{u}, \mathbf{w}) = F(\mathbf{u}) - F(\mathbf{w}) + (\mathbf{w} - \mathbf{u}, \nabla F(\mathbf{w})). \quad (4)$$

Here (\mathbf{u}, \mathbf{w}) denotes the usual dot product in \mathbb{R}^n .

The link function f is simply the gradient of F , i.e. $f = \nabla F$. Note that adding the constant term to F does not change the corresponding link function. Also, when a linear term is added to F , the link function changes by an additive constant. However, adding a linear term to F does not change the divergence D_F and the corresponding update (3).

The convex function $\text{sq}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^n w_i^2$ corresponds to the identity link, and the divergence for this choice of F is the squared Euclidean distance, denoted as $D_{\text{sq}}(\mathbf{u}, \mathbf{w}) = \sum_{i=1}^n \frac{1}{2} (u_i - w_i)^2$. Similarly, the convex function $\text{ne}(\mathbf{w}) = \sum_{i=1}^n w_i \ln w_i - w_i$ (essentially the negative entropy) gives rise to the logarithm link and the un-normalized relative entropy as the divergence function (see Figure 1): $D_{\text{ne}} = \sum_{i=1}^n (u_i \ln \frac{u_i}{w_i} + w_i - u_i)$.

The divergence functions D_F are used as potential functions for proving static loss bounds for the corresponding generalized gradient descent algorithms [CBLW96, KW97a, HKW95, Byl97, KW97b]. As in the previous work in convex programming, we use projections based on these divergences.

We now outline how worst-case loss bounds are obtained in the static case. At the center of all the static proofs [KW94, HKW95, Byl97, KW97b] for the generalized gradient descent algorithms lies the following type of inequality:

$$a(\eta) L_t(\mathbf{w}_t) - b(\eta) L_t(\mathbf{u}) \leq D_F(\mathbf{u}, \mathbf{w}_t) - D_F(\mathbf{u}, \mathbf{w}_{t+1}). \quad (5)$$

Here a and b are non-negative functions that depend on the learning rate η and upper bounds on the norms of the instances. Note that $D_F(\mathbf{u}, \mathbf{w}_t) - D_F(\mathbf{u}, \mathbf{w}_{t+1})$ may be seen as the progress towards the off-line weight vector \mathbf{u} . Since this inequality holds for each trial, we can sum Equation (5) over trials. The progresses towards \mathbf{u} that appear on the right-hand sides form a telescoping sum and only the first and last terms survive:

$$a(\eta) L(A, S) - b(\eta) L(\langle \mathbf{u}_1, \dots, \mathbf{u}_\ell \rangle, S) \leq D_F(\mathbf{u}, \mathbf{w}_1) - D_F(\mathbf{u}, \mathbf{w}_{\ell+1}). \quad (6)$$

The last term $D_F(\mathbf{u}, \mathbf{w}_{\ell+1})$ can be dropped, since the divergence is non-negative and we are forming an upper bound. The first term $D_F(\mathbf{u}, \mathbf{w}_1)$ plays the role of the size measure of \mathbf{u} . Rearranging the above and picking a good choice for η gives bounds of the form (1).

The new approach (which allows us to prove bounds for the shifting case) starts with picking a suitable convex region Γ . After the usual update the hypothesis vector is projected into this convex region. This gives us two benefits: first, we can show that the generalized gradient descent algorithms may be modified to maintain convex constraints on the hypothesis vector. Second, we show that shifting bounds may readily be obtained. The convex region plays a key role in bounding the size of the schedule of predictors.

Definition and properties of the Projection Update

The projection of a point x onto a closed convex set Γ is simply the point in Γ closest to x w.r.t. the divergence D_F .

Definition 1.1. *The projection of a point \mathbf{w} w.r.t. divergence $D_F : E \times \text{ri } E \rightarrow [0, \infty)$ where $E \subseteq \mathbb{R}^n$ onto a closed convex set Γ is defined by:*

$$P_{(\Gamma, D_F)}(\mathbf{w}) = \text{argmin}_{\mathbf{u} \in \Gamma \cap E} D_F(\mathbf{u}, \mathbf{w}). \quad (7)$$

In the Appendix we show that the projection exists and is unique.

Recall the standard Pythagorean Theorem for the divergence associated with the identity link:

$$\|u - w\|^2 = \|u - p\|^2 + \|p - w\|^2, \text{ when } u - p \perp p - w.$$

The orthogonality condition can be rewritten as p being a projection of w onto a plane that contains u . Here projections are w.r.t. the (squared) Euclidean distance. Surprisingly, the Pythagorean Theorem holds for projections w.r.t. D_F divergences.

Theorem 1.2. (Generalized Pythagoras) *Given a divergence $D_F : E \times ri E \rightarrow [0, \infty)$, a closed convex set Γ such that $\Gamma \cap ri E \neq \emptyset$, and points $w \in ri E$ and $u \in \Gamma$, then*

$$D_F(u, w) \geq D_F(u, P_{(\Gamma, D_F)}(w)) + D_F(P_{(\Gamma, D_F)}(w), w). \quad (8)$$

In the special case when Γ is a hyperplane then the above becomes an equality.

Note that this inequality is opposite to the triangle inequality that holds for a (distance) metric. This is the reason that we call D_F a divergence instead of a distance.

The above theorem has been proven many times under a variety of assumptions, see for instance [Bre67, Csi91, JB90]. For the sake of completeness we provide a streamlined proof in the Appendix. The below corollary of the Generalized Pythagorean Theorem will be used repeatedly in the analyses of our algorithms.

Corollary 1.3. *Given a divergence $D_F : E \times ri E \rightarrow [0, \infty)$, a closed convex set Γ such that $\Gamma \cap ri E \neq \emptyset$, and points $w \in ri E$ and $u \in \Gamma$, then*

$$D_F(u, w) - D_F(u, P_{(\Gamma, D_F)}(w)) \geq 0 \quad (9)$$

Using the above definition of projection we now introduce our modification of the generalized gradient descent algorithms:

$$\begin{aligned} w_t^m &= f^{-1}(f(w_t) - \eta \nabla L_t(w_t)) && \text{gen.grad.desc.upd.} \\ w_{t+1} &= P_{(\Gamma, D_F)}(w_t^m) && \text{proj.upd.} \end{aligned} \quad (10)$$

Since we now have two updates, we will refer to the intermediate weight vector following the generalized gradient update as w_t^m . We call this algorithm the Constrained Generalized Gradient Descent Algorithm. We use C-GD(Γ) and C-EGU(Γ) for the example algorithms of this paper.

The constraint sets Γ_{gd} and Γ_{egu} chosen for to obtain shifting bounds for GD and EGU, respectively, are intimately related to the static loss bounds proven for the original algorithms. In [KW97a] it is observed that the bounds of GD depend on the 2-norm of both the instances x_t and the predictor u , while in EGU (as well as EG) the bounds depend on the ∞ -norm of the instances x_t and the 1-norm of the predictor u . The constraint set is associated with the predictor and the norm dual to the predictor's norm. Thus the constraint set Γ_{gd} is the level set of the two-norm, and the constraint set Γ_{egu} is a shifted level set of the ∞ -norm. Our worst-case loss bounds for C-GD(Γ_{gd}) thus depend on

$\|\langle u_1, \dots, u_\ell \rangle\|_2$ and $\|x_t\|_2$, and the bounds for C-EGU(Γ_{egu}) depend on $\|\langle u_1, \dots, u_\ell \rangle\|_1$ of the predictor schedules and $\|x_t\|_\infty$.

Calculating the projection $P_{(\Gamma, D_F)}(w)$ may be computationally expensive. For our example algorithms the projections $P_{(\Gamma_{\text{gd}}, D_{\text{sq}})}(w)$ and $P_{(\Gamma_{\text{egu}}, D_{\text{ne}})}(w)$ are easily computable (see Figure 1) in $O(n)$ time. In particular, the projection update $P_{(\Gamma_{\text{egu}}, D_{\text{ne}})}(w)$ corresponds to simple weight clipping. Similar updates were used in the Fixed-share Algorithm [HW98], WML [LW94] and the algorithms for tracking disjunctions [AW98].

An overview of proof techniques based on projections

We have three results. First, we show that if the predictor u from the comparison class lies in the constraint set, then the previous static loss bounds hold for the constraint algorithm. Second, we prove a shifting bound on a predictor schedule $\langle u_1, \dots, u_\ell \rangle$ whose predictors must all lie in Γ . These results follow almost immediately from properties of projections and the static bounds. Third (more tricky), we show for C-EGU(Γ_{egu}) that the shifting bound may be extended from predictors in the constraint set $\Gamma_{\text{egu}} = [\frac{\gamma}{n}, \frac{n}{\gamma}]^n$ to predictors in the larger set $[0, \frac{n}{\gamma}]^n$. Interestingly, in this case the shifting loss bound obtained via the projection update of this paper is similar to the loss bound of the Fixed-share Algorithm [HW98].

For the first part recall that w_t is the weight vector at the start of a trial, and w_t^m the weight vector between the two updates. With this notation, Inequality (5) becomes

$$a(\eta)L_t(w_t) - b(\eta)L_t(u_t) \leq D_F(u, w_t) - D_F(u, w_t^m). \quad (15)$$

Note that this inequality no longer telescopes. By Corollary 1.3 we have

$$0 \leq D_F(u, w_t^m) - D_F(u, w_{t+1}),$$

provided that u lies in the constraint set. The sum of the two inequalities is again a telescoping inequality with the same functions a and b :

$$a(\eta)L_t(w_t) - b(\eta)L_t(u_t) \leq D_F(u, w_t) - D_F(u, w_{t+1}). \quad (16)$$

Thus the constraint algorithms have the same static bounds, provided that the comparison vector u lies in the constraint set.

For the second part we assume that all the predictors of predictor schedule $\langle u_1, \dots, u_\ell \rangle$ lie in the constraint set. Thus Equation (16) holds for each trial:

$$a(\eta)L_t(w_t) - b(\eta)L_t(u_t) \leq D_F(u_t, w_t) - D_F(u_t, w_{t+1}). \quad (17)$$

This again does not telescope. We fix this by adding

$$G_F(u_t, u_{t+1}) := D_F(u_t, w_{t+1}) - D_F(u_{t+1}, w_{t+1}). \quad (18)$$

Intuitively, $G_F(u_t, u_{t+1})$ measures the cost for shifting from u_t to u_{t+1} . We will return to a discussion of this cost later.

Convex Function $F(\mathbf{w})$	$\text{ne}(\mathbf{w}) = \sum_{i=1}^n w_i \ln w_i - w_i, \mathbf{w} \in [0, \infty)^n$	$\text{sq}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^n w_i^2, \mathbf{w} \in \mathbb{R}^n$
Link Function $f(\mathbf{w})$	$\ln(\mathbf{w}) = (\ln w_1, \dots, \ln w_n), \mathbf{w} \in (0, \infty)^n$	$\text{id}(\mathbf{w}) = \mathbf{w}, \mathbf{w} \in \mathbb{R}^n$
$D_F(\mathbf{u}, \mathbf{w})$	$D_{\text{ne}}(\mathbf{u}, \mathbf{w}) = \sum_{i=1}^n (u_i \ln \frac{u_i}{w_i} + w_i - u_i)$	$D_{\text{sq}}(\mathbf{u}, \mathbf{w}) = \sum_{i=1}^n \frac{1}{2} (u_i - w_i)^2$
Gen. Grad. Desc. upd.	$w_{t+1,i}^m = w_{t,i} e^{-\eta \nabla_{t,i}}$ (11)	$w_{t+1,i}^m = w_{t,i} - \nabla_{t,i}$ (12)
Constraint Set Γ	$\Gamma_{\text{egu}} = [\frac{\gamma}{n}, \frac{n}{\gamma}]^n$	$\Gamma_{\text{gd}} = \{\mathbf{u} : \ \mathbf{u}\ _2 \leq \gamma\}$
Projection Update	$\forall i: 1, \dots, n : w_{t+1,i} = \begin{cases} w_{t,i}^m & w_{t,i}^m \in [\frac{\gamma}{n}, \frac{n}{\gamma}] \\ \frac{\gamma}{n} & w_{t,i}^m < \frac{\gamma}{n} \\ \frac{n}{\gamma} & w_{t,i}^m > \frac{n}{\gamma} \end{cases}$ (13)	$w_{t+1} = \begin{cases} w_t^m & w_t^m \in \Gamma_{\text{gd}} \\ \frac{\gamma w_t^m}{\ w_t^m\ _2} & w_t^m \notin \Gamma_{\text{gd}} \end{cases}$ (14)

Figure 1: Two constrained generalized gradient descent algorithms C-EGU(Γ_{egu}) and C-GD(Γ_{gd}), based on the Unnormalized Exponentiated Gradient and the Gradient Descent. Here $\nabla_{t,i}$ is shorthand for $\frac{\partial L_t(\mathbf{w}_t)}{\partial w_{t,i}}$. In the case of the square loss, $\nabla_{t,i} = 2x_{t,i}(w_t \cdot \mathbf{x}_t - y_t)$.

By adding the Equation (17) and (18) and by summing over all trials we get the following:

$$L(A, S) \leq \frac{b(\eta)}{a(\eta)} L(\langle \mathbf{u}_1, \dots, \mathbf{u}_\ell \rangle, S) \left[\frac{1}{a(\eta)} [D_F(\mathbf{u}_1, \mathbf{w}_1) - D_F(\mathbf{u}_{\ell+1}, \mathbf{w}_{\ell+1}) - \sum_{t=1}^{\ell} G_F(\mathbf{u}_t, \mathbf{u}_{t+1})] \right]. \quad (19)$$

This would be a good bound if $G_F(\mathbf{u}_t, \mathbf{u}_{t+1})$ was lower bounded. Unfortunately this is not true for arbitrary \mathbf{w}_{t+1} . For example, for the squared Euclidean distance,

$$G_{\text{sq}}(\mathbf{u}_t, \mathbf{u}_{t+1}) = \|\mathbf{u}_t\|_2^2 - \|\mathbf{u}_{t+1}\|_2^2 - 2\mathbf{w}_{t+1} \cdot (\mathbf{u}_t - \mathbf{u}_{t+1}). \quad (20)$$

Observe that unless we have a bound on some norm of \mathbf{w}_{t+1} , the cost $G_{\text{sq}}(\mathbf{u}_t, \mathbf{u}_{t+1})$ cannot be lower bounded solely in terms of \mathbf{u}_t and \mathbf{u}_{t+1} . However, $\mathbf{w}_{t+1} \in \Gamma_{\text{gd}}$ and thus by our choice of the constraint set Γ_{gd} , $\|\mathbf{w}_{t+1}\|_2 \leq \gamma$. Through an application of Hölder's inequality, we have

$$G_{\text{sq}}(\mathbf{u}_t, \mathbf{u}_{t+1}) \geq \|\mathbf{u}_t\|_2^2 - \|\mathbf{u}_{t+1}\|_2^2 - 2\gamma \|\mathbf{u}_t - \mathbf{u}_{t+1}\|_2, \quad (21)$$

and plugging into Equation (19) we obtain a reasonable upper bound of the form (2):

$$L(\text{C-GD}(\Gamma_{\text{gd}}), S) \leq \frac{b(\eta)}{a(\eta)} L(\langle \mathbf{u}_1, \dots, \mathbf{u}_\ell \rangle, S) + \frac{1}{a(\eta)} [D_{\text{sq}}(\mathbf{u}_1, \mathbf{w}_1) - D_{\text{sq}}(\mathbf{u}_{\ell+1}, \mathbf{w}_{\ell+1}) + \|\mathbf{u}_{\ell+1}\|_2^2 - \|\mathbf{u}_1\|_2^2 + 2\gamma \sum_{t=1}^{\ell} \|\mathbf{u}_t - \mathbf{u}_{t+1}\|_2].$$

Part two for the divergence function D_{ne} is similar. Again the constraint set is used to obtain a lower bound of the G function: $G_{\text{egu}}(\mathbf{u}_t, \mathbf{u}_{t+1}) \geq -H(\mathbf{u}_t) + H(\mathbf{u}_{t+1}) - \|\mathbf{u}_t\|_1 +$

$\|\mathbf{u}_{t+1}\|_1 - \ln(\frac{n}{\gamma}) \|\mathbf{u}_t - \mathbf{u}_{t+1}\|_1$, where H denotes the entropy function. By plugging into (19) we obtain

$$L(\text{C-EGU}(\Gamma_{\text{egu}}), S) \leq \frac{b(\eta)}{a(\eta)} L(\langle \mathbf{u}_1, \dots, \mathbf{u}_\ell \rangle, S) + \frac{1}{a(\eta)} [D_{\text{ne}}(\mathbf{u}_1, \mathbf{w}_1) - D_{\text{ne}}(\mathbf{u}_{\ell+1}, \mathbf{w}_{\ell+1}) + H(\mathbf{u}_1) - H(\mathbf{u}_{\ell+1}) + \|\mathbf{u}_0\|_1 - \|\mathbf{u}_{\ell+1}\|_1 + \ln(\frac{n}{\gamma}) \sum_{t=1}^{\ell} \|\mathbf{u}_t - \mathbf{u}_{t+1}\|_1]. \quad (22)$$

The third part is more involved and is covered in detail in the rest of the paper. We only give a rough outline in this section. The analysis applies directly to EGU and extends the comparison class from $\Gamma_{\text{egu}} = [\frac{\gamma}{n}, \frac{n}{\gamma}]^n$ to $[0, \frac{n}{\gamma}]^n$. Since the comparison vectors do not lie in the constraint set we cannot use the Pythagorean Theorem. Instead, we directly lower bound $G_{\text{ne}}(\mathbf{u}_t, \mathbf{u}_{t+1})$. This lower bound will contain a constant term $-\gamma$, having the effect that our loss bounds for the extended concept class grow with the sequence length. Careful tunings of γ ameliorate this effect.

2 SHIFTING C-EGU(Γ_{egu}) FOR PREDICTORS IN $[0, \frac{n}{\gamma}]^n$

We begin by proving a lower bound for $G_{\text{ne}}(\mathbf{u}_t, \mathbf{u}_{t+1})$ as defined in Equation (18) when \mathbf{u}_t and \mathbf{u}_{t+1} lie in $[0, \frac{n}{\gamma}]^n$ instead of $\Gamma_{\text{egu}} = [\frac{\gamma}{n}, \frac{n}{\gamma}]^n$

Lemma 2.1. *Let $\mathbf{w} \in [0, \infty)^n$, and let $\mathbf{u}_t, \mathbf{u}_{t+1} \in [0, \frac{n}{\gamma}]^n$. Then*

$$G_{\text{ne}}(\mathbf{u}_t, \mathbf{u}_{t+1}) \geq -H(\mathbf{u}_t) + H(\mathbf{u}_{t+1}) - \|\mathbf{u}_t\|_1 + \|\mathbf{u}_{t+1}\|_1 - \ln(\frac{n}{\gamma}) \|\mathbf{u}_t - \mathbf{u}_{t+1}\|_1 - \gamma. \quad (23)$$

Proof. Note that the divergence D_{ne} is a simple sum over coordinates, and the constraint set Γ_{egu} is the same in each

coordinate. By a one-dimensional argument one can show that

$$P_{(\Gamma_{\text{egu}}, D_{\text{ne}})}(\mathbf{w}) = P_{([\frac{\gamma}{n}, \infty)^n, D_{\text{ne}})}(P_{([0, \frac{n}{\gamma}]^n, D_{\text{ne}})}(\mathbf{w})).$$

Let $\mathbf{w}' = P_{([0, \frac{n}{\gamma}]^n, D_{\text{ne}})}(\mathbf{w})$, and let $\mathbf{w}'' = P_{([\frac{\gamma}{n}, \infty)^n, D_{\text{ne}})}(\mathbf{w}')$. The sum of the following three inequalities gives Inequality (23), thus proving the lemma:

$$D_{\text{ne}}(\mathbf{u}_t, \mathbf{w}) - D_{\text{ne}}(\mathbf{u}_t, \mathbf{w}') \geq 0 \quad (24)$$

$$D_{\text{ne}}(\mathbf{u}_t, \mathbf{w}') - D_{\text{ne}}(\mathbf{u}_t, \mathbf{w}'') \geq -\gamma \quad (25)$$

$$D_{\text{ne}}(\mathbf{u}_t, \mathbf{w}'') - D_{\text{ne}}(\mathbf{u}_{t+1}, \mathbf{w}'') \geq -H(\mathbf{u}_t) + H(\mathbf{u}_{t+1}) - \|\mathbf{u}_t\|_1 + \|\mathbf{u}_{t+1}\|_1 - \ln\left(\frac{n}{\gamma}\right)\|\mathbf{u}_t - \mathbf{u}_{t+1}\|_1. \quad (26)$$

Inequality (24) holds by Corollary 1.3 when all \mathbf{u}_t lie in the constraint set $[0, \frac{n}{\gamma}]^n$. Recall that the latter convex set was used as the constraint set for defining \mathbf{w}' as a projection.

For Inequality (25) we expand the definition of D_{ne} :

$$D_{\text{ne}}(\mathbf{u}_t, \mathbf{w}') - D_{\text{ne}}(\mathbf{u}_t, \mathbf{w}'') = \sum_{i=1}^n (w'_i - w''_i) + \sum_{i=1}^n u_{t,i} \ln \frac{w'_i}{w''_i}$$

From the definition of \mathbf{w}' and \mathbf{w}'' it follows that $w'_i - w''_i \geq -\frac{\gamma}{n}$ and $w''_i - w'_i \geq 0$. Thus the first sum is lower bounded by $-\gamma$ and the second sum is lower bounded by zero, giving us Inequality (25).

For the proof of Inequality (26) we expand D_{ne} , apply Hölder's inequality and the fact that $w''_i \in [\frac{\gamma}{n}, \frac{n}{\gamma}]$:

$$\begin{aligned} D_{\text{ne}}(\mathbf{u}_t, \mathbf{w}'') - D_{\text{ne}}(\mathbf{u}_{t+1}, \mathbf{w}'') &= -H(\mathbf{u}_t) + H(\mathbf{u}_{t+1}) - \|\mathbf{u}_t\|_1 + \|\mathbf{u}_{t+1}\|_1 + \\ &\sum_{i=1}^n \left[(u_{t,i} - u_{t+1,i}) \ln \frac{1}{w''_i} \right] \\ &\geq -H(\mathbf{u}_t) + H(\mathbf{u}_{t+1}) - \|\mathbf{u}_t\|_1 + \|\mathbf{u}_{t+1}\|_1 - \\ &\ln\left(\frac{n}{\gamma}\right)\|\mathbf{u}_t - \mathbf{u}_{t+1}\|_1. \end{aligned}$$

□

By plugging the lower bound of the lemma into Equation (19) we obtain an upper bound which is the same as Equation (22), except for the additional term $\frac{\gamma\ell}{a(\eta)}$:

For all predictor schedules $\langle \mathbf{u}_1, \dots, \mathbf{u}_\ell \rangle \in \{[0, \frac{n}{\gamma}]^n\}^\ell$,

$$\begin{aligned} L(\text{C-EGU}(\Gamma_{\text{egu}}), S) &\leq \frac{b(\eta)}{a(\eta)} L(\langle \mathbf{u}_1, \dots, \mathbf{u}_\ell \rangle, S) \\ &+ \frac{1}{a(\eta)} [D_{\text{ne}}(\mathbf{u}_1, \mathbf{w}_1) - D_{\text{ne}}(\mathbf{u}_{\ell+1}, \mathbf{w}_{\ell+1}) + H(\mathbf{u}_1) - \\ &H(\mathbf{u}_{\ell+1}) + \|\mathbf{u}_0\|_1 - \|\mathbf{u}_{\ell+1}\|_1 + \ln\left(\frac{n}{\gamma}\right) \sum_{t=1}^{\ell} \|\mathbf{u}_t - \mathbf{u}_{t+1}\|_1 + \gamma\ell]. \end{aligned}$$

Finally, we give a tuning of γ for the purpose of minimizing the above. The tuning involves two parameters, $\hat{\ell}$ and \hat{U} , where $\hat{\ell}$ is an upper bound on the length of the trial sequence and \hat{U} is an approximation to $\|\mathbf{u}_1\|_1 + \|\langle \mathbf{u}_1, \dots, \mathbf{u}_\ell \rangle\|_1$. We set $\gamma = \frac{\hat{U}}{\hat{\ell} + \hat{U}}$ and $\mathbf{w}_1 = (\frac{\gamma}{n}, \dots, \frac{\gamma}{n})$. Then for all predictor

schedules $\langle \mathbf{u}_1, \dots, \mathbf{u}_\ell \rangle \in \{[0, \frac{n}{\gamma}]^n\}^\ell$, provided that $\hat{U} \leq \hat{\ell}$, we have the bound

$$\begin{aligned} L(\text{C-EGU}(\Gamma_{\text{egu}}), S) &\leq \frac{b(\eta)}{a(\eta)} L(\langle \mathbf{u}_1, \dots, \mathbf{u}_\ell \rangle, S) + \\ &\frac{1}{a(\eta)} (\|\mathbf{u}_1\|_1 + \|\langle \mathbf{u}_1, \dots, \mathbf{u}_\ell \rangle\|_1) (\ln n + \ln\left(\frac{\hat{\ell}}{\hat{U}}\right) + \ln 2) + \hat{U}. \end{aligned}$$

3 OPEN PROBLEMS

We feel that we are developing important methods based on projections that can be used for proving worst-case loss bounds when the predictor from the comparison class is allowed to shift. These methods apply to such algorithms as the WM Algorithm[LW94], the Aggregating Algorithm[Vov95], the Hedge Algorithm[FS97], and various exponentiated gradient algorithms [KW97a, HKW95, By197, KW97b], as well as Winnow [Lit88]. Just the static bounds for the constraint algorithms are already interesting in their own right. A number of detailed applications will be included in the full paper.

There were two algorithms for the shifting expert setting. The Fixed-share Algorithm, which essentially lower bounded the weights, and the Variable-share Algorithm, which lets the weights of the algorithm go to zero. So far we have only developed the methodology for the updates corresponding to the Fixed-share Algorithm (i.e. we use clipping to define the convex region Γ_{egu} associated with the exponentiated gradient algorithms). For the alternate bounds (in progress) we want to let the weights go close to zero and only pull back the weights when the master incurs large loss. Thus we need to use a dynamic convex region that we project into at end of each trial.

A nice application for the projection update may be made to prediction portfolios on-line[Cov91, HSSW96]. In this case linear inequality constraints may be used to express relations that must be maintained between the instruments of the portfolio.

Acknowledgments We would like to thank Claudio Gentile, Andrew Klingler, Debra Lewis and Harold Widom for valuable discussions.

References

- [AW98] P. Auer and M. K. Warmuth. Tracking the best disjunction. *Journal of Machine Learning*, 1998. Special issue on concept drift.
- [BB97] A. Blum and C. Burch. On-line learning and the metrical task system. In *Proc. 10th Annu. Workshop on Comput. Learning Theory*. ACM Press, New York, NY, 1997.
- [Bre67] L.M. Bregman. The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR Computational Mathematics and Physics*, 7:200–217, 1967.
- [By197] T. Bylander. The binary exponentiated gradient algorithm for learning linear functions. In *Proc. 8th Annu. Conf. on Comput. Learning Theory*, pages 184–192. ACM, 1997.

- [CBLW96] N. Cesa-Bianchi, P. Long, and M.K. Warmuth. Worst-case quadratic loss bounds for on-line prediction of linear functions by gradient descent. *IEEE Transactions on Neural Networks*, 7(2):604–619, May 1996.
- [CL81] Y. Censor and A. Lent. An iterative row-action method for interval convex programming. *Journal of Optimization Theory and Applications*, 34(3):321–353, July 1981.
- [Cov91] T. M. Cover. Universal portfolios. *Mathematical Finance*, 1(1):1–29, 1991.
- [Csi91] Imre Csiszar. Why least squares and maximum entropy? An axiomatic approach for linear inverse problems. *The Annals of Statistics*, 19(4):2032–2066, 1991.
- [FS97] Yoav Freund and Robert E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, August 1997.
- [GLS97] A. J. Grove, N. Littlestone, and D. Schuurmans. General convergence results for linear discriminant updates. In *Proc. 8th Annu. Conf. on Comput. Learning Theory*, pages 171–183. ACM, 1997.
- [HKW95] D. P. Helmbold, J. Kivinen, and M. K. Warmuth. Worst-case loss bounds for sigmoided linear neurons. In *Proc. 1995 Neural Information Processing Conference*, pages 309–315. MIT Press, Cambridge, MA, November 1995.
- [HKW97] D. Haussler, J. Kivinen, and M. K. Warmuth. Tight worst-case loss bounds for predicting with expert advice. *IEEE Transactions on Information Theory*, 1997. To appear.
- [HLS96] D.P. Helmbold, D.D.E. Long, and B. Sherrod. A dynamic disk spin-down technique for mobile computing. In *Proceedings of the Second Annual ACM International Conference on Mobile Computing and Networking*. ACM/IEEE, November 1996.
- [HSSW96] D. Helmbold, R. E. Schapire, Y. Singer, and M. K. Warmuth. On-line portfolio selection using multiplicative updates. In *Proc. 13th International Conference on Machine Learning*, pages 243–251. Morgan Kaufmann, San Francisco, July 1996.
- [HW98] M. Herbster and M. K. Warmuth. Tracking the best expert. *Journal of Machine Learning*, 1998. Special issue on concept drift.
- [JB90] L. Jones and C. Byrne. General entropy criteria for inverse problems, with applications to data compression, pattern classification and cluster analysis. *IEEE Transactions on Information Theory*, 36(1):23–30, 1990.
- [JW98] A. Jagota and M. K. Warmuth. Continuous and discrete time nonlinear gradient descent: relative loss bounds and convergence. In R. Greiner E. Boros, editor, *Electronic Proceedings of Fifth International Symposium on Artificial Intelligence and Mathematics*, pages – Electronic, <http://rutcor.rutgers.edu/~amai>, 1998.
- [KW94] J. Kivinen and M. Warmuth. Using experts for predicting continuous outcomes. In *Computational Learning Theory: Eurocolt '93*, volume New Series Number 53 of *The Institute of Mathematics and its Applications Conference Series*, pages 109–120, Oxford, 1994. Oxford University Press.
- [KW97a] J. Kivinen and M. K. Warmuth. Additive versus exponentiated gradient updates for linear prediction. *Information and Computation*, 132(1):1–64, January 1997.
- [KW97b] J. Kivinen and M. K. Warmuth. Relative loss bounds for multidimensional regression problems. In *Advances in Neural Information Processing Systems, 11*, Cambridge, MA, 1997. MIT Press. To appear.
- [Lit88] N. Littlestone. Learning when irrelevant attributes abound: A new linear-threshold algorithm. *Machine Learning*, 2:285–318, 1988.
- [Lit89] N. Littlestone. *Mistake Bounds and Logarithmic Linear-threshold Learning Algorithms*. PhD thesis, Technical Report UCSC-CRL-89-11, University of California Santa Cruz, 1989.
- [LW94] N. Littlestone and M. K. Warmuth. The weighted majority algorithm. *Information and Computation*, 108(2):212–261, 1994.
- [Roc70] R. Rockafellar. *Convex Analysis*. Princeton University Press, 1970.
- [Rud91] W. Rudin. *Functional Analysis*. McGraw-Hill, 1991.
- [Vov90] V. Vovk. Aggregating strategies. In *Proc. 3rd Annu. Workshop on Comput. Learning Theory*, pages 371–383. Morgan Kaufmann, 1990.
- [Vov95] V. G. Vovk. A game of prediction with expert advice. In *Proc. 8th Annu. Conf. on Comput. Learning Theory*, pages 51–60. ACM Press, New York, NY, 1995.
- [Vov97] V. Vovk. Derandomizing stochastic prediction strategies. In *Proc. 10th Annu. Workshop on Comput. Learning Theory*. ACM Press, New York, NY, 1997.

APPENDIX: DIVERGENCES AND THE GENERALIZED PYTHAGOREAN THEOREM

We begin by giving a more detailed definition of Bregman divergences. Here ri and bd denote the relative interior and relative boundary of a set [Roc70].

Definition A.1. Let F be a function from a closed convex set $E \subset \mathbb{R}^n$ into \mathbb{R} such that the following three conditions hold:

(C1) F is strictly convex.

(C2) F is differentiable.

(C3) $\forall r \in bd E, \forall s, t \in ri E : (\nabla F(r) - \nabla F(s), t - r) < 0$.

Then the Bregman divergence $D_F : E \times \text{ri } E \rightarrow [0, \infty)$ is defined as

$$D_F(\mathbf{u}, \mathbf{w}) = F(\mathbf{u}) - F(\mathbf{w}) + (\mathbf{w} - \mathbf{u}, \nabla F(\mathbf{w})). \quad (27)$$

The conditions that F is strictly convex and F is differentiable guarantee that $D_F(\mathbf{u}, \mathbf{w}) \geq 0$ and $\mathbf{u} = \mathbf{w} \Leftrightarrow D_F(\mathbf{u}, \mathbf{w}) = 0$.

The technical Condition (C3) guarantees that any projection (Definition 1.1) w.r.t. to the divergence will not lie on the boundary of E . Note that the inner product in Condition (C3) is the directional derivative of $D_F(\cdot, \mathbf{s})$ at the point \mathbf{r} in the direction $\mathbf{t} - \mathbf{r}$. The direction $\mathbf{t} - \mathbf{r}$ is the direction from \mathbf{r} (any point in $\text{bd } E$) to \mathbf{t} (any point in $\text{ri } E$).

The divergence $D_F(\mathbf{u}, \mathbf{w})$ is strictly convex in \mathbf{u} since it is the sum of a strictly convex function and a linear function. Similarly, $D_F(\mathbf{u}, \mathbf{w})$ is continuous in \mathbf{u} . We define the open and closed balls relative to a point and a divergence as $B_\epsilon(\mathbf{w}) = \{\mathbf{v} : D_F(\mathbf{v}, \mathbf{w}) < \epsilon\}$ and $\bar{B}_\epsilon(\mathbf{w}) = \{\mathbf{v} : D_F(\mathbf{v}, \mathbf{w}) \leq \epsilon\}$. Since $D_F(\mathbf{u}, \mathbf{w})$ is continuous in \mathbf{u} , the closure $B_\epsilon(\mathbf{w})$ is equal to $\bar{B}_\epsilon(\mathbf{w})$. It is straightforward to check that these balls are strictly convex sets since $D_F(\cdot, \mathbf{w})$ is strictly convex. In the full paper we will also show the following:

Proposition A.2. *For any point $\mathbf{w} \in \text{ri } E$, $\bar{B}_\epsilon(\mathbf{w})$ is bounded.*

The uniqueness of the projection essentially follows from the boundedness and strict convexity of the balls, and the convexity of the constraint set.

Proposition A.3. *Given $\mathbf{w} \in \text{ri } E$ and a closed convex set Γ such that $\Gamma \cap \text{ri } E \neq \emptyset$, then $P_{(\Gamma, D_F)}(\mathbf{w})$ exists and is unique.*

Proof. Without loss of generality, assume $\Gamma \subseteq E$. We first show existence. Let $A = D_F(\Gamma, \mathbf{w})$. Since all points in A are lower bounded by 0, $a = \inf A$ exists. We now prove a lies in A , which proves the existence of a projection. Let $\delta > 0$. The ball $\bar{B}_{a+\delta}(\mathbf{w})$ is closed and bounded. Therefore the intersection $B = \bar{B}_{a+\delta}(\mathbf{w}) \cap \Gamma$ is closed and bounded. Since we are in \mathfrak{R}^n , this implies that B is compact. Let $C = D_F(B, \mathbf{w})$. Since B is compact and $D_F(\cdot, \mathbf{w})$ is continuous, it follows that C is compact. Since C is compact, $\inf C$ must lie in C . It follows that a projection exists and all projections lie in $D = \{\mathbf{v} : D_F(\mathbf{v}, \mathbf{w}) = \inf C\} \cap \Gamma$. We now show that the set D has a single element, and hence the projection is unique. Let \mathbf{s}, \mathbf{t} be distinct elements of D and let $\mathbf{u} = \frac{\mathbf{s} + \mathbf{t}}{2}$. Then $\mathbf{u} \in \Gamma$, since Γ is convex. Since $D_F(\cdot, \mathbf{w})$ is strictly convex, $D_F(\mathbf{u}, \mathbf{w}) < \frac{D_F(\mathbf{s}, \mathbf{w})}{2} + \frac{D_F(\mathbf{t}, \mathbf{w})}{2} = \inf C$, which is a contradiction. Hence $P_{(\Gamma, D_F)}(\mathbf{w})$ is unique. \square

We prove a sequence of three propositions. The first and third comprise Theorem 1.2, the generalization of the Pythagorean Theorem. The following proposition treats the special case of this theorem when Γ is a hyperplane (second part of Theorem 1.2).

Proposition A.4. *Given a divergence $D_F : E \times \text{ri } E \rightarrow [0, \infty)$, a hyperplane $H = \{\mathbf{v} : (\mathbf{x}, \mathbf{v}) = y\}$ such that $H \cap \text{ri } E \neq \emptyset$, and points $\mathbf{w} \in \text{ri } E$ and $\mathbf{u} \in H$, then*

$$D_F(\mathbf{u}, \mathbf{w}) = D_F(\mathbf{u}, P_{(H, D_F)}(\mathbf{w})) + D_F(P_{(H, D_F)}(\mathbf{w}), \mathbf{w}). \quad (28)$$

Proof. If $\mathbf{w} \in H$, then $P_{(H, D_F)}(\mathbf{w}) = \mathbf{w}$ and the proposition trivially holds. Assume $\mathbf{w} \notin H$. By expanding the D_F (cf. Definition A.1) we get the following equivalent form of Equation (28):

$$(\mathbf{u} - P_{(H, D_F)}(\mathbf{w}), \nabla F(\mathbf{w}) - \nabla F(P_{(H, D_F)}(\mathbf{w}))) = 0. \quad (29)$$

We now proceed to prove an implicit form of the projection. Computing the projection is a convex programming problem, i.e. the computation of $\text{argmin}_{\mathbf{v}} D_F(\mathbf{v}, \mathbf{w})$ subject to the constraint that \mathbf{v} lies in the convex set $E \cap H$. We may ignore the constraint $\mathbf{v} \in E$, since technical Condition (C3) ensures that \mathbf{v} is in $\text{ri } E$. We introduce a Lagrange multiplier λ for the equality constraint that $\mathbf{v} \in H$:

$$L(\lambda, \mathbf{v}) = D_F(\mathbf{v}, \mathbf{w}) + \lambda((\mathbf{x}, \mathbf{v}) - y). \quad (30)$$

Thus for the projection the following equation holds:

$$\begin{aligned} 0 &= \nabla_{\mathbf{v}} L(\lambda, \mathbf{v}) \mid_{\mathbf{v} = P_{(H, D_F)}(\mathbf{w})} \\ &= \nabla F(P_{(H, D_F)}(\mathbf{w})) - \nabla F(\mathbf{w}) + \lambda \mathbf{x}. \end{aligned}$$

Rewriting the above equation, we have

$$\nabla F(\mathbf{w}) - \nabla F(P_{(H, D_F)}(\mathbf{w})) = \lambda \mathbf{x}, \quad (31)$$

and the left-hand-side of (29) becomes $\lambda(\mathbf{u} - P_{(H, D_F)}(\mathbf{w}), \mathbf{x})$. Since $(\mathbf{x}, P_{(H, D_F)}(\mathbf{w})) = y$ and $(\mathbf{x}, \mathbf{u}) = y$, this inner product is zero. \square

We next show the generalized Pythagorean Theorem for projections onto halfspaces.

Proposition A.5. *Given a divergence $D_F : E \times \text{ri } E \rightarrow [0, \infty)$, a closed halfspace $\mathcal{H} = \{\mathbf{v} : (\mathbf{x}, \mathbf{v}) \leq y\}$ such that $\mathcal{H} \cap \text{ri } E \neq \emptyset$, and points $\mathbf{w} \in \text{ri } E$ and $\mathbf{u} \in \mathcal{H}$, then*

$$D_F(\mathbf{u}, \mathbf{w}) \geq D_F(\mathbf{u}, P_{(\mathcal{H}, D_F)}(\mathbf{w})) + D_F(P_{(\mathcal{H}, D_F)}(\mathbf{w}), \mathbf{w}). \quad (32)$$

Proof. The case $\mathbf{w} \in \mathcal{H}$ is trivially true. Assume $\mathbf{w} \notin \mathcal{H}$. Equation (32) has the following equivalent form:

$$(\mathbf{u} - P_{(\mathcal{H}, D_F)}(\mathbf{w}), \nabla F(\mathbf{w}) - \nabla F(P_{(\mathcal{H}, D_F)}(\mathbf{w}))) \leq 0. \quad (33)$$

The projection of \mathbf{w} clearly lies on the boundary of the halfspace, i.e. the hyperplane $H = \{\mathbf{v} : (\mathbf{x}, \mathbf{v}) = y\}$. Thus $P_{(\mathcal{H}, D_F)}(\mathbf{w}) = P_{(H, D_F)}(\mathbf{w})$. Finding the projection is a convex programming problem subject to an inequality constraint. Thus the parameter λ in (30) is now a Kuhn-Tucker coefficient and Equation (31) again holds. This lets us rewrite the left-hand-side of (33) as $\lambda(\mathbf{u} - P_{(H, D_F)}(\mathbf{w}), \mathbf{x})$. We complete the proof by showing that this expression is at most zero. First note that the projection lies on the boundary H of \mathcal{H} and thus the Kuhn-Tucker coefficient is positive. Finally, $(\mathbf{x}, P_{(H, D_F)}(\mathbf{w})) = y$ and $(\mathbf{x}, \mathbf{u}) \leq y$, so we have that

$$\lambda(\mathbf{u} - P_{(H, D_F)}(\mathbf{w}), \mathbf{x}) \leq 0. \quad \square$$

We now prove the first part of the generalized Pythagorean Theorem:

Proposition A.6. *Given a divergence $D_F : E \times \text{ri } E \rightarrow [0, \infty)$, a closed convex set Γ such that $\Gamma \cap \text{ri } E \neq \emptyset$, and points $\mathbf{w} \in \text{ri } E$ and $\mathbf{u} \in \Gamma$, then*

$$D_F(\mathbf{u}, \mathbf{w}) \geq D_F(\mathbf{u}, P_{(\Gamma, D_F)}(\mathbf{w})) + D_F(P_{(\Gamma, D_F)}(\mathbf{w}), \mathbf{w}). \quad (34)$$

Proof. The case $\mathbf{w} \in \Gamma$ is again trivially true. Assume $\mathbf{w} \notin \Gamma$. Let $\epsilon = D_F(P_{(\Gamma, D_F)}(\mathbf{w}), \mathbf{w})$. Consider the balls $B_\epsilon(\mathbf{w})$ and $\overline{B}_\epsilon(\mathbf{w})$. Note that $B_\epsilon(\mathbf{w}) \cap \Gamma = \emptyset$ and $\overline{B}_\epsilon(\mathbf{w}) \cap \Gamma = P_{(\Gamma, D_F)}(\mathbf{w})$. Let \mathcal{H} be a halfspace that contains Γ but does not intersect with $B_\epsilon(\mathbf{w})$. The existence of such a halfspace is implied by the Hahn-Banach separation Theorem ([Rud91, p59]). (This theorem says that for any disjoint open and closed sets there exists a halfspace containing the closed set such that this halfspace is disjoint with the open set.)

We now claim that $P_{(\Gamma, D_F)}(\mathbf{w}) = P_{(\mathcal{H}, D_F)}(\mathbf{w})$. Provided this claim is true, the current proposition clearly follows from the previous one. Since $\mathcal{H} \supseteq \Gamma$,

$$D_F(P_{(\mathcal{H}, D_F)}(\mathbf{w}), \mathbf{w}) \leq D_F(P_{(\Gamma, D_F)}(\mathbf{w}), \mathbf{w}) = \epsilon.$$

On the other hand $<$ is impossible because $B_\epsilon(\mathbf{w})$ and \mathcal{H} are disjoint. Thus the above inequality must be an equality and claim follows from the uniqueness of projections. \square