



UvA-DARE (Digital Academic Repository)

Combining strategies efficiently: high-quality decisions from conflicting advice

Koolen, W.M.

Publication date
2011

[Link to publication](#)

Citation for published version (APA):

Koolen, W. M. (2011). *Combining strategies efficiently: high-quality decisions from conflicting advice*. [Thesis, fully internal, Universiteit van Amsterdam].

General rights

It is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), other than for strictly personal, individual use, unless the work is under an open content license (like Creative Commons).

Disclaimer/Complaints regulations

If you believe that digital publication of certain material infringes any of your rights or (privacy) interests, please let the Library know, stating your reasons. In case of a legitimate complaint, the Library will make the material inaccessible and/or remove it from the website. Please Ask the Library: <https://uba.uva.nl/en/contact>, or a letter to: Library of the University of Amsterdam, Secretariat, P.O. Box 19185, 1000 GD Amsterdam, The Netherlands. You will be contacted as soon as possible.

$$\operatorname{argmin}_{w \in \operatorname{conv}(\mathcal{C})} \Delta(w \| w^{t-1}) + \eta w \cdot \ell^t$$

Abstract We develop an online algorithm called *Component Hedge* for learning structured concept classes when the loss of a structured concept sums over its components. Example classes include paths through a graph (composed of edges) and partial permutations (composed of assignments). The algorithm maintains a parameter vector with one non-negative weight per component, which always lies in the convex hull of the structured concept class. The algorithm predicts by decomposing the current parameter vector into a convex combination of concepts and choosing one of those concepts at random. The parameters are updated by first performing a multiplicative update and then projecting back into the convex hull. We show that Component Hedge has optimal regret bounds for a large variety of structured concept classes.

6.1 Introduction

We develop online learning algorithms for structured concepts that are composed of components. For example, sets are composed of elements, permutations of individual assignments, trees have edges as components, etc. The number of components d is considered small, but the number of structured concepts D built from the components is typically exponential in d .

Our algorithms address the following online prediction problem. In each trial the algorithm first produces a concept from the structured class by choosing a concept probabilistically based on its current parameters. It then observes the loss of each concept. Finally, it prepares for the next trial by updating its parameters by incorporating the losses. Since the algorithm “hedges” by choosing the structured concept probabilistically, we analyse the expected loss incurred in each trial. The goal is to develop algorithms with small regret, which is the total expected loss of the online algorithm minus the loss of the best structured concept in the class chosen in hindsight.

We now make a key simplifying assumption on the loss: We assume that the loss of a structured concept in each trial is always the sum of the losses of its components and that the component losses always have range $[0, 1]$. Thus if the concepts are k -element sets chosen out of n elements, then in each trial each element is assigned a loss in $[0, 1]$ and the loss of any particular k -set is simply the sum of the losses of its elements. Similarly for trees, a loss in $[0, 1]$ is assigned to each edge of the graph and the loss of a tree is the sum of the losses of its edges.

We will show that with this simplifying assumption we still have rich learning problems that address a variety of new settings. We give efficient algorithms (i.e. polynomial in d) that serve as an entry point for considering more complex losses in the future.

Perhaps the simplest approach to learn structured concept classes online is the Follow the Perturbed Leader (FPL) algorithm [92]. FPL adds a random perturbation to the cumulative loss of each individual component, and then plays the structured concept with minimal perturbed loss. FPL is widely applicable, since efficient combinatorial optimisation algorithms exist for a broad range of concept classes. Unfortunately, the loss range of the structured concepts enters into the regret bounds that we can prove for FPL. For example, for k -sets the

loss range is $[0, k]$ because each set contains k elements, for permutations the loss range is $[0, n]$ because each permutation is composed of n assignments, etc.

A second simple approach for learning well compared to the best structured concept is to run the Hedge algorithm of [59] with one weight per structured concept. The original algorithm was developed for the so-called expert setting, which in the context of this chapter corresponds to learning with sets of size one. To apply this algorithm to our setting, the experts are chosen as the structured concepts in the class we are trying to learn. In this chapter we call this algorithm *Expanded Hedge* (EH). It maintains its uncertainty as a probability distribution over all structured concepts and the weight W_C of concept C is proportional to $\exp(-\eta\ell(C))$, where $\ell(C)$ is the total loss of concept C incurred so far and η is a non-negative learning rate.

There are two problems with EH. First, there are exponentially many weights to maintain. However our simplifying assumption assures that $\ell(C)$ is a sum over the losses of the component of C . This implies that W_C is proportional to a product over the components of the structured concept C and this fact can be exploited to still achieve efficient algorithms in some cases. More importantly however, like for FPL, the loss range of the structured concepts usually enters into the best regret bounds that we can prove.

Learning with structured concepts has also been dealt with recently in the bandit domain [26]. However all of this work is based on EH and contains the additional range factors.

Our contribution Our new method, called *Component Hedge* (CH), avoids the additional range factors altogether. Each structured concept C is identified with its incidence vector in $\{0, 1\}^d$ indicating which components are used. The parameter space of CH is simply the convex hull of all concepts in the class \mathcal{C} to be learned. Thus, whereas EH maintains a weight for each structured concept, CH only maintains a weight for each component. The current parameter vector represents CH's first-order "uncertainty" about the quality of each concept. The value of parameter i represents the *usage* of component i in the next prediction. The usages of the components are updated in each trial by incorporating the current losses, and if the usage vector leaves the hull, then it is projected back via a relative entropy projection. The key trick

to make this projection efficient is to find a representation of the convex hull of the concepts as a convex polytope with a number of facets that is polynomial in d . We give many applications where this is possible.

We clearly champion the Component Hedge algorithm in this chapter because we can prove regret bounds for this algorithm that are tight within constant factors for many structured concept classes. Also it is trivial to enhance CH with a variety of “share updates” that make it robust in the case when the best comparator changes over time [80, 19].

Two instances of CH have appeared before even though this name was not used: learning with k -sets [185] and learning with permutations [77]. The same polytope we use for paths was also employed in [5] for developing online algorithms for the bandit setting. They avoid the projection step altogether by exploiting a barrier function. The contribution of this chapter is to clearly formulate the general methodology of the Component Hedge algorithm and give many more involved combinatorial examples. In the case of permutations we also show how the method can be used to learn truncated permutations. Also in earlier work [173] it was pointed out that the Expanded Hedge algorithm can be simulated efficiently in many cases. In particular, the concept class of paths in a directed graph was introduced. However, good bounds were only achieved in very special cases. In this chapter we show that CH essentially is optimal for the path problem.

Outline We give the basic setup for the structured prediction task, introduce CH and prove its general regret bound in Section 6.2. We then turn to a list of applications in Section 6.3: vanilla experts, k -sets, permutations, paths, undirected and directed spanning trees. For each structured concept class we discuss efficient implementation of CH, and derive expected regret bounds for this algorithm. Then in Section 6.4 we provide matching lower bounds for all examples, showing that the regret of CH is optimal within a constant factor. In Section 6.5 we compare CH to the existing algorithms EH and FPL. We observe that the best general regret bounds for each algorithm exceed that of CH by a significant range factor. We show that the bounds for these other algorithms can be improved to closely match those of CH whenever the so-called *unit rule* holds for the algorithms and class. This means any loss vector $\ell \in [0, 1]^d$ can be split into up to d scaled unit loss vectors $\ell_i e_i$ and processing these in separate trials always incurs at least

as much loss. Unfortunately, for most pairing of the algorithms CH and FPL with the classes we consider in this chapter, we have explicit counter examples to the unit rule. Finally, Section 6.6 concludes with a list of open problems.

6.2 Component Hedge

Prediction task We consider sequential prediction [75, 25] based on a structured concept class [92, 26]. Fix a set of concepts $\mathcal{C} \subseteq \{0, 1\}^d$ of size $D = |\mathcal{C}|$. For example \mathcal{C} could consist of the incidence vectors of subsets of k out of n elements (then $D = \binom{n}{k}$ and $d = n$), or the adjacency matrices of undirected spanning trees on n elements (then $D = n^{n-2}$ and $d = n(n-1)/2$).

Our online learning protocol proceeds in trials. At trial t , we have to produce a single concept $C^t \in \mathcal{C}$. Then a loss vector $\ell^t \in [0, 1]^d$ is revealed, and we incur loss given by the dot product $C^t \cdot \ell^t$. Although each component suffers loss at most 1, a concept may suffer loss up to $U := \max_{C \in \mathcal{C}} |C|$. We allow randomised algorithms. Thus the expected loss of the algorithm at trial t is $\mathbb{E}[C^t] \cdot \ell^t$, where the expectation is over the internal randomisation of the algorithm. Our goal is to minimise our (expected) *regret* after T trials

$$\sum_{t=1}^T \mathbb{E}[C^t] \cdot \ell^t - \min_{C \in \mathcal{C}} \sum_{t=1}^T C \cdot \ell^t.$$

That is, the difference between our cumulative expected loss and the loss of the best concept in hindsight.

Note that the i th component of $\mathbb{E}[C^t]$ is the probability that component i is “used in” concept C^t . We therefore call $\mathbb{E}[C^t]$ the *usage vector*. This vector becomes the internal parameter of our algorithm. The set of all usage vectors is the convex hull of the concepts.

6.2.1 Component Hedge

Two instances of CH appeared before in the literature [77, 185]. Here we give the algorithm in its general form, and prove a general regret bound. The algorithm CH maintains its uncertainty about the best structured concept as a usage vector w^t in $\text{conv}(\mathcal{C}) \subseteq [0, 1]^d$, the convex

Table 6.1 Example structured concept classes

Case	U	D	d
Experts	1	n	n
k -Sets	k	$\binom{n}{k}$	n
Permutations	n	$n!$	n^2
Paths (from source to sink)	$n + 1$	$n! \cdot e - o(1)$	$n(n + 1) + 1$
Undirected spanning trees	$n - 1$	n^{n-2}	$n(n - 1)/2$
Directed spanning trees	$n - 1$	n^{n-2}	$(n - 1)^2$

hull of the concepts \mathcal{C} . The initial weight w^0 is typically the usage of the uniform distribution on concepts. CH predicts in trial t by decomposing w^{t-1} into a convex combination¹ of the concepts \mathcal{C} , then sampling \mathbf{C}^t according to its weight in that convex combination. The expected loss of CH is thus $w^{t-1} \cdot \ell^t$. The updated weight w^t is obtained by trading off the relative entropy with the linear loss:

$$w^t := \operatorname{argmin}_{w \in \operatorname{conv}(\mathcal{C})} \Delta(w \| w^{t-1}) + \eta w \cdot \ell^t,$$

where the relative entropy is defined by

$$\Delta(w \| v) = \sum_{i \in [d]} \left(w_i \ln \frac{w_i}{v_i} + v_i - w_i \right).$$

It is easy to see that this update can be split into two steps: an unconstrained update followed by relative entropy projection into the convex hull:

$$\begin{aligned} \hat{w}^t &:= \operatorname{argmin}_{w \in \mathbb{R}^d} \Delta(w \| w^{t-1}) + \eta w \cdot \ell^t \\ w^t &:= \operatorname{argmin}_{w \in \operatorname{conv}(\mathcal{C})} \Delta(w \| \hat{w}^t). \end{aligned}$$

It is easy to see that $\hat{w}_i^t = w_i^{t-1} e^{-\eta \ell_i^t}$, that is, the old weights are simply scaled down by the exponentiated losses. The result of the relative

¹This decomposition usually is far from unique.

entropy projection w^t unfortunately does not have a closed form expression.

For CH to be efficiently implementable, the hull has to be captured by polynomial in d many constraints. This will allow us to efficiently decompose any point in the hull as a convex combination of at most $d + 1$ concepts. The trickier part is to efficiently implement the projection step. For this purpose one can use generic convex optimisation routines. For example this was done in the context of implementing the entropy regularised boosting algorithm [186]. We proceed on a case by case basis and often develop iterative algorithms that locally enforce constraints and do multiple passes over all constraints. See Table 6.1 for a list of structured concept classes we consider in this chapter.

6.2.2 Regret Bounds

As in [77], the analysis is split into two steps parallelling the two update steps. Essentially the unnormalised update step already gives the regret bound and the projection step does not hurt. For any usage vector $w^{t-1} \in \text{conv}(\mathcal{C})$, loss vector $\ell^t \in \{0, 1\}^d$ and any comparator concept C ,

$$\begin{aligned} (1 - e^{-\eta})w^{t-1} \cdot \ell^t &\leq \underbrace{\Delta(C\|w^{t-1}) - \Delta(C\|\hat{w}^t)}_{\sum_i w_i^{t-1}(1 - e^{-\eta \ell_i^t})} + \eta C \cdot \ell^t \\ &\leq \Delta(C\|w^{t-1}) - \Delta(C\|w^t) + \eta C \cdot \ell^t \end{aligned}$$

The first inequality is obtained by bounding the exponential using the inequality $1 - e^{-\eta x} \geq (1 - e^{-\eta})x$ for $x \in [0, 1]$ as done in [108]. The second inequality is an application of the Generalised Pythagorean Theorem [81], using the fact that w^t is a Bregman projection of \hat{w}^t into the convex set $\text{conv}(\mathcal{C})$, which contains C . We now sum over trials and obtain, abbreviating $\ell^1 + \dots + \ell^T$ to $\ell^{\leq T}$,

$$(1 - e^{-\eta}) \sum_{t=1}^T w^{t-1} \cdot \ell^t \leq \Delta(C\|w^0) - \Delta(C\|w^T) + \eta C \cdot \ell^{\leq T}.$$

Recall that $w^{t-1} \cdot \ell^t$ equals the expected loss $\mathbb{E}[C^t] \cdot \ell^t$ of CH in trial t . Also, relative entropies are nonnegative, so we may drop the second one, giving us the following bound on the total loss of the algorithm:

$$\sum_{t=1}^T \mathbb{E}[C^t] \cdot \ell^t \leq \frac{\Delta(C\|w^0) + \eta C \cdot \ell^{\leq T}}{1 - e^{-\eta}}.$$

To proceed we have to expand the prior w^0 . We consider the *symmetric balanced* case, i.e. where the concept class is invariant under permutation of the components, and every concept uses exactly U components. Paths may have different lengths and hence do not satisfy these requirements. All other examples from Table 6.1 do. In this balanced symmetric case we take w^0 to be the usage of the uniform distribution on concepts, satisfying $w_i^0 = U/d$ for each component i . It follows that $\Delta(C\|w^0) = U \ln(d/U)$, because any comparator C is a 0/1 vector that also uses exactly U components.

Let ℓ^* denote $\min_{C \in \mathcal{C}} C \cdot \ell^{\leq T}$, the loss of the best concept in hindsight. Then by choosing $\eta = \sqrt{\frac{2U \ln(d/U)}{\ell^*}}$ as a function of ℓ^* , we obtain the following general expected regret bound for CH:

$$\mathbb{E}[\ell_{\text{CH}}] - \ell^* \leq \sqrt{2\ell^*U \ln(d/U)} + U \ln(d/U). \quad (6.1)$$

The best-known general regret bounds for Expanded Hedge [59] and Follow the Perturbed Leader [84] are:

$$\mathbb{E}[\ell_{\text{EH}}] - \ell^* \leq \sqrt{2\ell^*U \ln D} + U \ln D \quad (6.2)$$

$$\mathbb{E}[\ell_{\text{FPL}}] - \ell^* \leq \sqrt{4\ell^*Ud \ln d} + 3Ud \ln d \quad (6.3)$$

where $D = |\mathcal{C}|$. Specific values for U , D and d in each application are listed in Table 6.1. We remark that if only an upper bound $\hat{\ell} \geq \ell^*$ is available, then we can still tune η as a function of $\hat{\ell}$ to achieve these bounds with $\hat{\ell}$ under the square roots instead of ℓ^* . Moreover, standard heuristics can be used to tune η “online” when no good upper bound on ℓ^* is given, which increase the expected regret bounds by at most a constant factor. (e.g. [28, 84]).

We are not concerned with small multiplicative constants (e.g. 2 vs 4), but the gap between (6.1) and both (6.2) and (6.3) is significant. To compare, observe that $\ln D$ is of order $U \ln d$ in all our applications. Thus, the EH regret bound is worse by a factor \sqrt{U} , while FPL is worse by a bigger factor \sqrt{d} . Moreover, in Section 6.4 we show for the covered examples that our expected regret bound (6.1) for CH is optimal up to constant scaling.

Some concept classes have special structure that can be exploited to improve the regret bounds of FPL and EH down to that of CH. We consider one such property, called the *unit rule* in Section 6.5.

6.3 Applications

We consider the following structured concept classes: experts, k -sets, truncated permutations, source-sink paths, and both undirected and directed spanning trees. In each case we discuss implementation of CH and obtain a regret bound. Matching lower bounds are presented in Section 6.4.

6.3.1 Experts

The most basic example is the vanilla expert setting. In this case, the set of “structured” concepts equals the set of n standard basis vectors in \mathbb{R}^n . We will see that in this case Component Hedge see gracefully degrades to the original Hedge algorithm. First, the parameter spaces of both algorithms coincide since the convex hull of the basis vectors equals the probability simplex. Second, the predictions coincide since a vector in the probability simplex decomposes uniquely into a convex combination of basis vectors. Third, the parameter updates are the same, since the relative entropy projection of a non-negative weight vector into the probability simplex amounts to re-normalising to unity.

In fact on this simple task CH, EH and FPL each coincide with Hedge. For CH and EH this is obvious. For FPL this fact was observed in [102, 91] by using log-of-exponential perturbations instead of exponential perturbations used in the original paper [92]. Thus, we obtain following regret bound for all algorithms:

$$\mathbb{E}[\ell_{\text{CH}}] - \ell^* \leq \sqrt{2\ell^* \ln n} + \ln n.$$

6.3.2 k -sets

The problem of learning with sets of k out of n elements was introduced in [185] and applied to online Principal Component Analysis (PCA). Their algorithm is an instance of CH, and we review it here. The convex hull of k -sets equals the set of $w \in \mathbb{R}_+^n$ that satisfy the following constraints:

$$w_i \leq 1 \quad \text{for all } i \in [n] \quad \text{and} \quad \sum_{i=1}^n w_i = k. \quad (6.4)$$

Relative entropy projection into this polytope amounts to renormalising the sum to k , followed by redistributing the mass of the components that exceed 1 over the remaining components so that their ratios are preserved. Finally, each element of the convex hull of sets can be greedily decomposed into a convex combination of n k -sets by iteratively removing sets in the convex combination while always setting the coefficient of the new set as high as possible. Both projection and decomposition take $O(n^2)$ time [185].

Regret bound By (6.1), the regret of CH on sets is

$$\mathbb{E}[\ell_{\text{CH}}] - \ell^* \leq \sqrt{2\ell^*k \ln(n/k)} + k \ln(n/k).$$

We give a matching lower bound in Section 6.4.

6.3.3 Truncated Permutations

The second instantiation of CH that has appeared is the problem of permutations [77]. Here we consider a slightly generalised task: *truncated permutations* of k out of n elements. A truncated permutation fills k slots with distinct elements from a pool of n elements. Equivalently, a truncated permutation is a maximal matching in the complete bipartite graph between $[k]$ and $[n]$. Truncated permutations extend k -sets by linearly ordering the selected k elements.

Results to search queries are usually in the form of a truncated permutation; of all n existing documents, only the top k are displayed in order of decreasing relevance. Predicting with truncated permutations is thus a model for learning the best search result.

Matching polytope We write $i \leftarrow j$ for the component that assigns item j to slot i . Now the convex hull of truncated permutations consists of all $w \in \mathbb{R}_+^{k \times n}$ (see [161, Corollary 18.1b]) satisfying the following k row (left) and n column (right) constraints:

$$\sum_{j \in [n]} w_{i \leftarrow j} = 1 \quad \text{for all } i \in [k] \quad \text{and} \quad \sum_{i \in [k]} w_{i \leftarrow j} \leq 1 \quad \text{for all } j \in [n]. \quad (6.5)$$

Relative entropy projection The relative entropy projection of \widehat{w} into the convex hull of truncated permutations $w = \operatorname{argmin}_{w \text{ s.t. (6.5)}} \Delta(w \parallel \widehat{w})$ has no closed form solution. By convex duality (details are given in Appendix 6.B.1), $w_{i \leftarrow j} = \widehat{w}_{i \leftarrow j} e^{-\lambda_i - \mu_j}$, where λ_i and μ_j are the Lagrange multipliers associated to the row and column constraints (6.5), which minimise

$$\sum_{i \in [k]; j \in [n]} \widehat{w}_{i \leftarrow j} e^{-\lambda_i - \mu_j} + \sum_{i \in [k]} \lambda_i + \sum_{j \in [n]} \mu_j.$$

under the constraint that $\mu \geq \mathbf{0}$. This dual problem, which has $2n$ variables and n constraints, may be optimised directly using numerical convex optimisation software. Another approach is to iteratively reestablish each violated constraint beginning from $\mu = \mathbf{0}$ and $\lambda = \mathbf{0}$. In full permutation case ($k = n$), this process is called *Sinkhorn balancing*. It is known to converge to the optimum, see [77] for an overview of efficiency and convergence results of this iterative method.

Decomposition Our decomposition algorithm for truncated permutations interpolates between the decomposition algorithms used for k -sets and full permutations [185, 77]. Assume w lies in the hull of truncated permutations, i.e. the constraints (6.5) are satisfied. To measure progress, we define a score $s(w)$ as the number of zero components in w plus the number of column constraints that are satisfied with equality.

Our algorithm maintains a truncated permutation C that satisfies the following invariant: C hits all columns whose constraints are satisfied with equality by w , and avoids all components with weight zero in w . Such a C can be established in time $O(k^2 n)$ using augmenting path methods (see [161, Theorem 16.3]).

Let l be the minimum weight of the components used by C , and let h be the maximum column sum of the columns untouched by C . So by construction $h < 1$. If $l = 1$ then $w = C$ and we are done. Otherwise, let $\alpha = \min\{l, 1 - h\}$, and set $w' = (w - \alpha C) / (1 - \alpha)$. It is easy to see that the vector w' satisfies (6.5), and that $s(w') > s(w)$. It is no longer the case that C satisfies the invariant w.r.t. w' . However, we may compute a weight k matching C' that satisfies the invariant by executing at most $s(w') - s(w)$ many augmenting path computations, which each cost $O(kn)$ time. We describe how this works below. After that we simply recurse on w' and C' . The resulting convex combination is αC plus $(1 - \alpha)$ times the result of the recursion.

The number of iterations is bounded by the score $s(w)$, which is at most kn . Thus, the total running time is $O(k^2n^2)$.

We now show that C can be improved to C' satisfying the invariant by a single augmenting path computation per violated requirement. Let C^* be a size k matching satisfying the invariant for w' . Such a matching always exists because w' lies in the matching polytope. Let $j \in [n]$ be a problematic column, i.e. either C matches j to a row i but $w'_{i \leftarrow j} = 0$, or C does not match j while its column constraint is tight for w' . From j , alternately follow edges from C and C^* . Since C and C^* are both matchings, this can not lead to a cycle, so it must lead to a path. Since all rows are matched, this path must end at a column. The path can not end at a column whose constraint is forced in both C and C^* . So it must end at a column whose constraint is not tight. Incorporating this augmenting path into C corrects the violated requirement without creating any new violations.

Regret bound By (6.1), the regret of CH on truncated permutations is

$$\mathbb{E}[\ell_{\text{CH}}] - \ell^* \leq \sqrt{2\ell^*k \ln n} + k \ln n.$$

We obtain a matching lower bound in Section 6.4.

6.3.4 Paths

The online shortest path problem was considered by [173, 92], and by various researchers in the bandit setting (see e.g. [26, 5] and references therein). We develop expected regret bounds for CH for the “full information setting”. Our regret bound improves the bounds given in [173, 92] which have the additional range factors in the square root.

Consider the a directed graph on the set of nodes $[n] \cup \{s, t\}$. Each trial we have to play a walk from the source node s to the sink node t . As always, our loss is given by the sum of the losses of the edges that our walk traverses. Since each edge loss is nonnegative (it lies in $[0, 1]$ by assumption) it is never beneficial to visit a node more than once. Thus w.l.o.g. we restrict attention to paths.

As an example, consider the full directed graph on $[n] \cup \{s, t\}$. Paths of length $k + 1$ through this graph use k distinct internal nodes in order, and therefore are in 1-1 correspondence with truncated permutations

of size k . Paths thus generalise truncated permutations by allowing all lengths simultaneously.

Unit flow polytope To implement CH efficiently, we have to succinctly describe the convex hull of paths. Unfortunately, we can not hope to write down linear constraints that capture the convex hull *exactly*. For if we could, then we could solve the *longest path* problem, which is known to be NP complete, by linear programming. Fortunately, there is a slight relaxation of the convex hull of paths that is describable by few constraints, namely the polytope of so-called unit flows. Even better, we will see that this relaxation does not hurt predictive performance at all.

A *unit flow* $w \in \mathbb{R}_+^d$ is described by the following constraints:

$$1 = \sum_{j \in [n]+t} w_{s,j} \quad \text{and} \quad \sum_{j \in [n]+s} w_{j,i} = \sum_{j \in [n]+t} w_{i,j} \quad \text{for each } i \in [n]. \quad (6.6)$$

We think of $w_{i,j}$ as describing the amount of flow from node i to j . The left constraint ensures that one unit of flow leaves the source s . The right constraint enforces that at internal nodes inflow equals outflow. It easily follows that one unit of flow enters the sink t .

The unit flow polytope is not bounded, but it has the right “bottom”. Namely, the vertices of the unit flow polytope are the s - t paths, see [161, Section 10.3]. The unit flow polytope is the Minkowski sum of the convex hull of s - t paths and the conic hull (nonnegative linear combinations) of directed cycles. Moreover, each unit flow can be decomposed into at most d paths and cycles, by iterative greedy removal of a directed cycle or paths containing the edge of least non-zero weight in time $O(n^4)$.

Since the unit flow polytope does have polynomially many constraints, we may efficiently run CH on it. Each round, it produces a flow. We then decompose this flow into paths and cycles, and throw away the cycles. We then sample a path from the remaining convex combination of paths.

Relative entropy projection To run CH, we have to compute the relative entropy projection of an arbitrary vector in \mathbb{R}_+^d into the flow polytope (6.6). This is a convex optimisation problem in $d \approx n^2$ variables

with constraints. By Slater's constraint condition, we have strong duality. So equivalently, we may solve the concave dual problem, which only has $n + 1$ variables and is unconstrained. The dual problem (details are given in Appendix 6.B.2) can therefore be solved efficiently by numerical convex optimisation software.

Say we want to find w , the relative entropy projection of \widehat{w} into the flow polytope. Since each edge appears in exactly two constraints with opposite sign, the solution has the form $w_{i,j} = \widehat{w}_{i,j} e^{\lambda_i - \lambda_j}$ for all $i, j \in [n] \cup \{s, t\}$, where λ_i is the Lagrange multiplier associated with node i (and $\lambda_t = 0$). The vector λ maximises

$$\lambda_s - \sum_{i \neq t; j \neq s} \widehat{w}_{i,j} e^{\lambda_i - \lambda_j}$$

That is, we have to find a single scale factor e^{λ_i} for each node i , such that scaling each edge weight by the ratio of the factors of its nodes reestablishes the flow constraints (6.6).

We propose the following iterative algorithm. Start with all λ_i equal to zero. Then pick a violated constraint, say at node i , and reestablish it by changing its associated λ_i . That is, we execute either

$$e^{\lambda_s} \leftarrow \frac{1}{\sum_{j \in [n] + t} \widehat{w}_{s,j} e^{-\lambda_j}}$$

or

$$e^{\lambda_i} \leftarrow \sqrt{\frac{\sum_{j \in [n] + s} \widehat{w}_{j,i} e^{\lambda_j}}{\sum_{j \in [n] + t} \widehat{w}_{i,j} e^{-\lambda_j}}} \quad \text{for some } i \in [n].$$

In our experiments, this algorithm converges quickly. We leave its thorough analysis as an open problem.

Decomposition Find any s - t path with non-zero weights on all edges in time $O(n^2)$. Subtract that path, scaled by its minimum edge weight. This creates a new zero, maintains flow balance, and reduces the outflow of the source. After at most n^2 iterations the source has outflow zero. Discard the remaining conic combination of directed cycles. The total running time is $O(n^4)$.

Regret bound for the complete directed graph Since paths have different lengths, we aim for a regret bound that depends on the length of the comparator path. To get such a bound, we need a prior usage vector w^0 that favours shorter paths. To this end, consider the distribution \mathbb{P} that distributes weight 2^{-k} uniformly over all paths of length $k \leq n$, and assigns weight 2^{-n} to the paths of length $n + 1$. This assures that \mathbb{P} is normalised to 1. Since there are $n!/(n - k + 1)!$ paths of length k , the probability of a path P of length k equals

$$\mathbb{P}(\mathbf{P} = P) = \begin{cases} \frac{(n - k + 1)!}{2^k n!} & \text{if } k \leq n, \\ \frac{1}{2^n n!} & \text{if } k = n + 1. \end{cases}$$

Also, the expected path length $\mathbb{E}[\mathbf{P} \cdot \mathbf{1}]$ is $2 - 2^{-n}$. We now set $w^0 := \mathbb{E}[\mathbf{P}]$, i.e. the usage of \mathbb{P} . There are three kinds of edges. We have one direct edge s, t , we have $2n$ boundary edges of the form s, j or i, t , and we have $n(n - 1)$ internal edges of the type i, j . A simple computation shows that their usages are (for $n \geq 3$)

$$w_{s,t}^0 = \frac{1}{2}, \quad w_{s,j}^0, w_{i,t}^0 = \frac{1}{2n}, \quad w_{i,j}^0 = \frac{1 - 2^{-(n-1)}}{2n(n-1)}.$$

Let P be a comparator path of length k . If $k = 1$ then $\Delta(P \| w^0) = \ln 2$. Otherwise, still for $n \geq 3$,

$$\begin{aligned} \Delta(P \| w^0) &= -2 \ln \frac{1}{2n} - (k - 2) \ln \frac{1 - 2^{-(n-1)}}{2n(n-1)} + \mathbb{E}[\mathbf{P} \cdot \mathbf{1}] - k \\ &= (k - 2) \ln(2n(n - 1)) + 2 \ln 2n + (k - 2) \ln \left(1 + \frac{2^{-(n-1)}}{1 - 2^{-(n-1)}} \right) - \\ &2^{-n} - (k - 2) \leq k \ln 2 - (k - 2) \frac{1 - 2^{-n+2}}{1 - 2^{-n+1}} + 2(k - 1) \ln n \leq 2k \ln n. \end{aligned}$$

By tuning η as before, the regret of CH with prior w^0 w.r.t. a comparator path of length k is

$$\mathbb{E}[\ell_{\text{CH}}] - \ell^* \leq \sqrt{4\ell^* k \ln n} + 2k \ln n.$$

This new regret bound improves known results in two ways. First, it does not have the range factors, which in the case of paths usually turn

out to be the diameter of the graph, i.e. the length of the longest s - t path. Second, some previous bounds only hold for acyclic graphs. Our bound holds for the complete graph.

Regret bound for an arbitrary graph We discussed the full graph as a first application of CH. For prediction on an arbitrary graphs we simply design a prior w^0 with zero usage on all edges that are not present in the graph. We could either use graph-specific knowledge, or we could use our old w^0 , disable edges by setting their usage to zero, and project back into the flow polytope. Relative entropy projection never revives zeroed edges. The regret bound now obviously depends on the graph via the prior usage w^0 .

6.3.4.1 Expanded Hedge and Component Hedge are Different on Paths

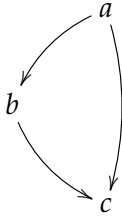
An efficient dynamic programming-based algorithm for EH was presented in [173]. This algorithm keeps one weight per edge, just like CH. These weights are updated using the *weight pushing algorithm*. This algorithm performs relative entropy projection on full distributions on paths. Like CH, weight pushing finds a weight of each node, and scales each edge weight by the ratio of its nodes weights. We now show that CH and EH are different on graphs. Consider the graph shown in Figure 6.1a. Say we use prior \mathbb{P} with weight $1/2$ on both paths (a, b, c) and (a, c) . Then the usages are $(1/2, 1/2, 1/2)$ for (ab, bc, ac) . Now multiply edge ab by $1/3$ (that is, we give it loss $\ln 3$), and both other edges by 1 (we give them loss zero). The resulting usages of EH and CH are displayed in Table 6.1b. The usages are different, and hence, so are the expected losses. In most cases (as shown e.g. in Table 6.1c), the updated usages of CH are irrational while the prior usages and the scale factors of the update are rational. On the other hand, EH always maintains rationality.

6.3.5 Spanning Trees

Whereas paths connect the source to the sink, spanning trees connect every node to every other node. Undirected spanning trees are often used in network-level communication protocols. For example, the *Spanning Tree Protocol* (IEEE 802.1D) is used by mesh networks of Ethernet

Figure 6.1 Expanded Hedge is not Component Hedge on paths

(a) Graph



(b) Usages after update (1/3, 1, 1)

Case	ab, bc	ac
EH and CH prior	1/2	1/2
EH after update	1/4	3/4
CH after update	1/3	2/3

(c) Usages after update (1/2, 1, 1)

Case	ab, bc	ac
EH and CH prior	1/2	1/2
EH after update	1/3	2/3
CH after update	$\frac{\sqrt{17}-1}{8}$	$\frac{9-\sqrt{17}}{8}$

switches to agree on a single undirected spanning tree, and thus eliminate loops by disabling redundant links. Directed spanning trees are used for asymmetric communication, for example for streaming multimedia from a central server to all connected clients. In either case, the cost of a spanning tree is the sum of the costs of its edges.

Learning spanning trees was pioneered by [99] for learning dependency parse trees. They discuss efficient methods for parameter estimation under log-loss and hinge loss. [26] derive a regret bound for undirected spanning trees in the bandit setting. We instantiate CH to both directed and undirected trees and give the first regret bound without the range factor.

Three kinds of directed spanning trees are common. Spanning trees with a fixed root, spanning trees with a single arbitrary root, and arborescences (or spanning forests) with multiple roots. We focus on a fixed root. The other two models can be simulated by a fixed root. To simulate arborescences, add a dummy as the fixed root, and put the root selection cost of node i along the path from the dummy to i . Furthermore, to force a single root, increase the cost of all edges leaving the dummy by a fixed huge amount.

Tree polytope To characterise the convex hull of directed trees on n nodes with fixed root 1, we use a trick based on flows from [112] that makes use of auxiliary variables $f_{i,j}^k$. For $i, j, k \in [n]$ the constraints are

$$0 \leq f_{i,j}^k \leq w_{i,j}, \quad \sum_{i,j} w_{i,j} = n - 1, \quad (6.7a)$$

and

$$\underbrace{\sum_{j \neq i} f_{j,i}^k}_{k\text{-flow into } i} + \underbrace{\mathbf{1}_{i=1}}_{k\text{-source at } 1} = \underbrace{\sum_{j \neq i} f_{i,j}^k}_{k\text{-flow out of } i} + \underbrace{\mathbf{1}_{i=k}}_{k\text{-sink at } k}. \quad (6.7b)$$

The intuition is as follows. A tree has $n - 1$ edges, and every node can be reached from the root. We enforce this by having a separate flow channel f^k for each non-root node k . We place a unit of flow into this channel at the root. Each intermediate node satisfies flow equilibrium. Finally, the target node k consumes the unit of flow destined for it. The first equation ensures that each edge's usage is sufficient for the flow that traverses that edge. The undirected tree polytope is constructed based on the directed tree polytope by considering the above $w_{i,j}$ as auxiliary variables, and imposing the constraint $w_{i,j} + w_{j,i} = v_{i,j}$. Now v are the weights sought.

Relative entropy projection The relative entropy projection of \widehat{w} into the convex hull of directed spanning trees $w = \operatorname{argmin}_{w \text{ s.t. (6.7)}} \Delta(w \| \widehat{w})$ has no closed form solution. By convex duality (details are given in Appendix 6.B.3), the solution satisfies

$$w_{i,j} = (n - 1) \frac{\widehat{w}_{i,j} e^{\sum_{k \neq 1} \max\{0, \mu_j^k - \mu_i^k\}}}{\sum_{i,j \neq i} \widehat{w}_{i,j} e^{\sum_{k \neq 1} \max\{0, \mu_j^k - \mu_i^k\}}}, \quad f_{ij}^k = \begin{cases} w_{i,j} & \text{if } \mu_j^k > \mu_i^k, \\ 0 & \text{if } \mu_j^k < \mu_i^k, \end{cases}$$

where μ_i^k , the Lagrange multipliers associated to the flow balance constraints, maximise

$$\sum_{k \neq 1} (\mu_k^k - \mu_1^k) - (n - 1) \ln \left(\sum_{i,j \neq i} \widehat{w}_{i,j} e^{\sum_{k \neq 1} \max\{0, \mu_j^k - \mu_i^k\}} \right).$$

This unconstrained concave maximisation problem in $\approx n^2$ variables seems easier than the primal problem, which has $\approx n^3$ variables and

constraints. Note however that the objective is not differentiable everywhere. Alternatively, we may again proceed by iteratively reestablishing constraints locally, starting from some initial assignment to the dual variables μ . This approach is analogous to Sinkhorn balancing.

Decomposition We have no special-purpose tree decomposition algorithm, and therefore resort to a general decomposition algorithm for convex polytopes that is based on linear programming. Let w be in the tree polytope. Choose an arbitrary vertex C (i.e. a spanning tree) by minimising a linear objective over the current polytope. Now use linear programming to find the furthest point w' in the polytope on the ray from C through w . At least one more inequality constraint is tight for w' . Thus w' lies in a convex polytope of at least one dimension lower. Add this inequality constraint as an equality constraint, recursively decompose w' , and express w as a convex combination of C and the decomposition of w' . The recursion bottoms out at a vertex (i.e. a spanning tree) and the total number of iterations is at most d .

Regret bound By (6.1), the regret $\mathbb{E}[\ell_{\text{CH}}] - \ell^*$ of CH on undirected and directed spanning trees is at most

$$\begin{aligned} \mathbb{E}[\ell_{\text{CH}}] - \ell^* &\leq \sqrt{2\ell^*(n-1)\ln(n/2) + (n-1)\ln(n/2)}, \\ \mathbb{E}[\ell_{\text{CH}}] - \ell^* &\leq \sqrt{2\ell^*(n-1)\ln(n-1) + (n-1)\ln(n-1)}. \end{aligned}$$

We provide matching lower bounds in Section 6.4.

6.4 Lower Bounds

Whereas it is easy to get some regret bounds with additional range factors, we show that CH is essentially optimal in all our applications. We leverage the following lower bound for the vanilla expert case:

6.4.1. THEOREM. *There are positive constants c_1 and c_2 s.t. any online algorithm for q experts with loss range $[0, U]$ can be forced to have expected regret at least*

$$c_1\sqrt{\ell^*U\ln q} + c_2\ln q. \tag{6.8}$$

This type of bound was recently proven in [3]. Note that c_1 and c_2 are independent of the number of experts, the range of the losses and the algorithm. Earlier versions of the above lower bound using many quantifier and limit arguments are given in [28, 77]. We now prove lower bounds for our structured concept classes by embedding the original expert problem into each class and applying the above theorem. This type of reduction was pioneered in [77] for permutations.

The general reduction works as follows. We identify q structured concepts C_1, \dots, C_q in the concept class $\mathcal{C} \subseteq \{0, 1\}^d$ to be learned that partition the d components. Now assume we have an online algorithm for learning class \mathcal{C} . From this we construct an algorithm for learning with q experts with loss range $[0, U]$. Let $\ell \in [0, U]^q$ denote the loss vector for the expert setting. From this we construct a loss vector $L \in [0, 1]^d$ for learning \mathcal{C} : $L := \sum_{i=1}^q \frac{\ell_i}{U} C_i$. That is, we spread the loss of expert i , evenly among the U many components used by concept C_i . Second, we transform the predictions as follows. Say our algorithm for learning \mathcal{C} predicts with any structured concept $C \in \mathcal{C}$. Then we play expert i with probability $C_i \cdot C/U$. The expected loss of the expert algorithm now equals the transformed loss of the algorithm for learning concepts in \mathcal{C} :

$$\mathbb{E}[\ell_i] = \sum_{i=1}^q \frac{C_i \cdot C}{U} \ell_i = C \cdot \sum_{i=1}^q \frac{\ell_i}{U} C_i = C \cdot L$$

This also means that the expected loss of the expert algorithm equals the expected loss of the algorithm for learning the structured class. This implies that the expected regret of the algorithm for learning \mathcal{C} is at least the expected regret of the expert algorithm. The lower bound (6.8) for the regret in the expert setting is thus also a lower bound for the regret of the structured prediction task.

k -sets We assume that k divides n . Then we can partition $[d]$ with n/k sets, where set i uses components $(i-1)k+1, \dots, ik$. The resulting lower bound has leading factor $\sqrt{k \ln \frac{n}{k}}$, matching the upper bound for CH within constant factors.

Truncated permutations We can partition the n^2 assignments into n full permutations. For example, the n cyclic shifts of the identity per-

mutation achieve this. The truncations to length k of those n permutations partition the kn components in the truncated case. The lower bound with leading factor $\sqrt{k \ln n}$ again matches the regret bound of CH within constant factors.

Spanning trees As observed in [72], the complete undirected graph has $(n-1)/2$ edge-disjoint spanning trees. Hence we get a lower bound with leading factor $\sqrt{(n-1) \ln((n-1)/2)}$. Each undirected spanning tree can be made directed by fixing a root. So there are at least as many disjoint directed spanning trees with a fixed root. In both cases we match the regret of CH within a constant factor.

Paths Consider the directed graph on $[n] \cup s, t$ that has n/k disjoint s - t paths of length $k+1$ connecting source to sink. By construction, we can embed n/k experts with loss range $[0, k]$ into this graph, so the regret has leading factor at least $\sqrt{k \log(n/k)}$. This graph is a subgraph of the complete directed graph $s \rightarrow K_n \rightarrow t$. Moreover, nature can force the algorithm to essentially play on the disjoint path graph by giving all edges outside it sheer infinite loss in a sheer infinite number of trials. This shows that the regret w.r.t. a comparator path of length k through the full graph has leading factor at least $\sqrt{k \log(n/k)}$.

A lower bound on the regret for arbitrary graphs is difficult to obtain since various interesting problems can be encoded as path problems. For example, the expert problem where each expert has a different loss range can be encoded into a graph that has a disjoint path of each length $1, 2, \dots, n$. The optimal algorithm for such expert problems was recently found in [2], but its regret has no closed form expression. It might be that the regret of CH is tight within constant factors for all graphs, but this question remains open.

6.5 Comparison to Other Algorithms

CH is a new member of an existing ecosystem. Other algorithms for structured prediction are EH[108] and FPL [92]. We now compare them.

Efficiency FPL can be readily applied efficiently to our examples of structured concept classes: k -sets take $O(n)$ per trial using variants

of median-finding, truncated permutations take $O(k^2n)$ per trial using the Hungarian method for minimum weight bipartite matching, paths take $O(n^2)$ per trial using Dijkstra's shortest path algorithm and spanning trees take $O(n^2)$ per trial using either Prim's algorithm or Chu–Liu/Edmonds's algorithm for finding a minimum weight spanning tree.

EH can be efficiently implemented for k -sets [185] and paths [173] using dynamic programming, and for spanning trees [99] using the Matrix-Tree Theorem by Kirchoff (undirected) and Tutte (directed). An approximate implementation based on MCMC sampling could be built for permutations based upon [85].

In most cases FPL and EH are faster than CH. This may be partly due to the novelty of CH and the lack of special-purpose algorithms for it. On the other hand, FPL solves a linear minimisation problem, which is intuitively simpler than minimising a convex relative entropy.

6.5.1 Improved Regret Bounds with the Unit Rule

On the other hand, we saw in Section 6.2.2 that the general regret bound for CH (6.1) improves the guarantees of EH (6.1) by a factor \sqrt{U} and those of FPL (6.3) by a larger factor \sqrt{d} . It is an open question whether these factors are real or simply an artifact of the bounding technique (see Section 6.6). We now give an example of a property of structured concept classes that makes these range factors vanish.

We say that a prediction algorithm has the *unit rule* on a given structured concept class \mathcal{C} if its worst-case performance is achieved when in each trial only a single expert has nonzero loss. Without changing the prediction algorithm, the unit rule immediately improves its regret bound by reducing the effective loss range of each concept from $[0, U]$ to $[0, 1]$. The improved regret bounds are (c.f. (6.2) and (6.3))

$$\mathbb{E}[\ell_{\text{EH}}] \leq \ell^* + \sqrt{2\ell^* \ln D} + \ln D \quad (6.9)$$

$$\mathbb{E}[\ell_{\text{FPL}}] \leq \ell^* + \sqrt{4\ell^* U \ln d} + 3U \ln d \quad (6.10)$$

The unit rules for EH and FPL on experts have been observed before [92, 6]. We reprove them here for completeness. The unit rule holds for both EH and FPL on sets, and for EH on undirected trees. It fails for EH and FPL on permutations, and for EH on directed trees.

We prove the unit rule for EH on sets here, and counter it for EH on directed trees. All other proofs and counterexamples are delayed to Appendix 6.A.

6.5.1.1 Unit Rule Holds for EH on k -sets

Fix an expert i , and let j be an arbitrary other expert. We claim that if we hand out loss to i , then the usage of j increases. For each k -set S , we denote the prior weight of S by W_S . We abbreviate

$$\begin{aligned} Z_i &:= \sum_{S:i \in S} W_S, & Z_{\neg i} &:= \sum_{S:i \notin S} W_S, \\ Z_j &:= \sum_{S:j \in S} W_S, & Z_{\neg j} &:= \sum_{S:j \notin S} W_S, \\ Z_{i \wedge j} &:= \sum_{S:i \in S, j \in S} W_S, & Z_{\neg i \wedge j} &:= \sum_{S:i \notin S, j \in S} W_S, \\ Z_{i \wedge \neg j} &:= \sum_{S:i \in S, j \notin S} W_S, & Z_{\neg i \wedge \neg j} &:= \sum_{S:i \notin S, j \notin S} W_S. \end{aligned}$$

6.5.1. THEOREM. *Assume that the prior weights have product structure, i.e. $W_S \propto \prod_{i \in S} w_i$. Then*

$$Z_j = \mathbb{P}(j \in \mathbf{S}^1) \leq \mathbb{P}(j \in \mathbf{S}^2 | \ell^1 = \delta_i) = \frac{Z_{i \wedge j} e^{-\eta} + Z_{\neg i \wedge j}}{Z_i e^{-\eta} + Z_{\neg i}}.$$

Proof. With some rewriting, the claim is equivalent to

$$Z_i Z_j \geq Z_{i \wedge j} \quad \text{and also} \quad Z_{i \wedge \neg j} Z_{\neg i \wedge j} \geq Z_{i \wedge j} Z_{\neg i \wedge \neg j}$$

Define

$$R(n, k) := \sum_{\substack{S \subseteq [n] \\ |S|=k}} \prod_{i \in S} w_i.$$

We now show that $R(n, k+1)R(n, m) \geq R(n, k)R(n, m+1)$ for all $0 \leq k < m < n$. The proof proceeds by induction on n . The case $n = 0$ is trivial. Now suppose that the claim holds up to n . We need to show it for $n+1$. For $n > 0$, we have

$$R(n, k) = \mathbf{1}_{k>0} w_n R(n-1, k-1) + \mathbf{1}_{k<n} R(n-1, k). \quad (6.11)$$

Suppose that the induction hypothesis holds up to n . We must show that for all $0 \leq k < m < n + 1$

$$R(n + 1, k + 1)R(n + 1, m) \geq R(n + 1, k)R(n + 1, m + 1).$$

By (6.11), this is equivalent to

$$\begin{aligned} & (w_{n+1}R(n, k) + \mathbf{1}_{k < n}R(n, k + 1))(\mathbf{1}_{m > 0}w_{n+1}R(n, m - 1) + \mathbf{1}_{m \leq n}R(n, m)) \geq \\ & (\mathbf{1}_{k > 0}w_{n+1}R(n, k - 1) + \mathbf{1}_{k \leq n}R(n, k))(\mathbf{1}_{m+1 > 0}w_{n+1}R(n, m) + \mathbf{1}_{m < n}R(n, m + 1)). \end{aligned}$$

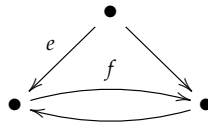
Now we expand, and use $0 \leq k < m < n + 1$ to eliminate indicators. It remains to show

$$\begin{pmatrix} (w_{n+1})^2R(n, k)R(n, m - 1) + \\ w_{n+1}R(n, k)R(n, m) + \\ w_{n+1}R(n, k + 1)R(n, m - 1) + \\ R(n, k + 1)R(n, m) \end{pmatrix} \geq \begin{pmatrix} \mathbf{1}_{k > 0}(w_{n+1})^2R(n, k - 1)R(n, m) + \\ \mathbf{1}_{k > 0}\mathbf{1}_{m < n}w_{n+1}R(n, k - 1)R(n, m + 1) + \\ w_{n+1}R(n, k)R(n, m) + \\ \mathbf{1}_{m < n}R(n, k)R(n, m + 1) \end{pmatrix}$$

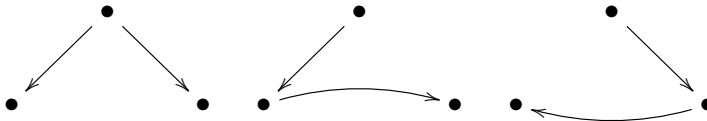
We now show that this inequality holds line-wise. Lines with active indicators trivially hold. If $k - 1 = m$, the second line holds with equality. Otherwise, and for the other lines we use the induction hypothesis. \square

6.5.1.2 Unit Rule Fails for EH on Directed Spanning Trees

The unit rule is violated for EH on directed trees. Consider this graph



and its three directed spanning trees:



Note that we may always restrict attention to a given graph G by assigning zero prior weight to all spanning trees of the full graph that use edges outside G . Now if we put a unit of loss on edge e , the usage of f decreases, and vice versa, contradicting the unit rule. Call the prior

weights on directed trees W_A, W_B, W_C . Then the usages satisfy

$$W_A + W_B = \mathbb{P}(e \in \mathbf{T}^1) \geq \mathbb{P}(e \in \mathbf{T}^2 | \ell^1 = \delta_f) = \frac{W_A + W_B e^{-\eta}}{W_A + W_B e^{-\eta} + W_C},$$

$$W_B = \mathbb{P}(f \in \mathbf{T}^1) \geq \mathbb{P}(f \in \mathbf{T}^2 | \ell^1 = \delta_e) = \frac{W_B e^{-\eta}}{W_A e^{-\eta} + W_B e^{-\eta} + W_C}.$$

Hence the unit rule is violated on directed spanning trees.

6.6 Conclusion

We developed the Component Hedge algorithm for online prediction over structured expert classes. The advantage of CH is that it has a general regret bound without the range factors that typically plague EH and FPL. We considered several example concept classes, and showed that the lower bound is matched in each case.

Open problems While the unit rule is one method for proving regret bounds for EH and FPL that are close to optimum, there might be other proof methods that show that EH and FPL perform as well as CH when applied to structured concepts. We know of no examples of structured concept classes where EH and FPL are clearly suboptimal. Resolving the question of whether such examples exist is our main open problem.

The prediction task for each structured concept class can be analysed as a two-player zero-sum game versus nature which tries to maximise the regret. The paper [6] gave an efficient implementation of the minimax optimal algorithm for playing against an adversary in the vanilla expert setting. Actually, the key insight was that the unit rule holds for the optimal algorithm in the vanilla expert case. This fact made it possible to design a balanced algorithm that incurs the same loss no matter which sequence of unit losses is chosen by nature. Unfortunately, the optimum algorithm does not satisfy the unit rule for any of the structured concept classes considered here. However, there might be some sort of relaxation of the unit rule that still leads to an efficient implementation of the optimum algorithm.

In this chapter the loss of a structured concept C always had the form $C \cdot \ell$, where ℓ is the loss vector for the components. This allowed us to maintain a mixture of concepts w and predict with a random

concept \mathbf{C} s.t. $\mathbb{E}[\mathbf{C}] = w$. By linearity, the expected loss of such a randomly drawn concept \mathbf{C} is the same as the loss of the mixture w . For regression problems with for example the convex loss $(C \cdot \ell - y)^2$ our algorithm can still maintain a mixture w , but now the expected loss of \mathbf{C} , i.e. $\mathbb{E}[(\mathbf{C} \cdot \ell - y)^2]$, is typically larger than the loss $(w \cdot \ell - y)^2$ of the mixture. We are confident that in this more general setting we can still get good regret bounds compared to the best mixture chosen in hind-sight. All we need to do is replace CH with the more general “Component Exponentiated Gradient” algorithm, which would do an EG update on the parameter vector w and project the updated vector back into the hull of the concepts.

In general, we believe that we have a versatile method of learning with structured concept classes. For example it is easy to augment the updates with a “share update” [80, 19] for the purpose of making them robust against sequences of examples where the best comparator changes over time. We also believe that our methods will get rid of the additional range factors in the bandit setting [26] and that gain versions of the algorithm CH also have good regret bounds.

At the core of our methods lies a relative entropy regularization which results in a multiplicative update on the components. In general, which relative entropy to choose is always one of the deepest questions. For example in the case of learning k -sets, a sum of binary relative entropies over the component can be used that incorporates the $w_i \leq 1$ constraints into the relative entropy term. In general incorporating inequality constraints into the relative entropy seems to have many advantages. However how to do this is an open ended research question.

6.A Unit rule

This appendix gives proofs of and counterexamples to the unit rule. We already saw some unit rule results in Section 6.5.1.

6.A.1 Unit Rule Holds for EH on Experts

Let \mathbf{E}^1 and \mathbf{E}^2 denote a random expert sampled by the algorithm in the first and second trial. Let δ_j denote a loss vector that assigns unit loss to expert j and zero loss to all other experts. Let W_i and W_j be the prior

weight of experts i and j . Then

$$W_j = \mathbb{P}(j = \mathbf{E}^1) \leq \mathbb{P}(j = \mathbf{E}^2 | \ell^1 = \delta_i) = \frac{W_j}{1 - W_i(1 - e^{-\eta})}.$$

Thus, if we hand out loss to one expert, *all* other usages increase. This unit rule result does not lead to improved regret bounds, (6.2) and (6.9) already coincide for experts.

6.A.2 Unit Rule Fails for EH on Permutations

There are two permutations of size two: $A = \{(1 \leftarrow 1), (2 \leftarrow 2)\}$ and $B = \{(1 \leftarrow 2), (2 \leftarrow 1)\}$. To contradict the unit rule, we show that if we give a unit of loss to the component $(1 \leftarrow 1)$, then the usage of $(2 \leftarrow 2)$ goes down with it and vice versa. By symmetry, we only need to show it in one order. Let W_A, W_B denote the prior weights of the two permutations. Then the usages satisfy

$$W_A = \mathbb{P}\left((2 \leftarrow 2) \in \Pi^1\right) \geq \mathbb{P}\left((2 \leftarrow 2) \in \Pi^1 \mid \ell^1 = \delta_{(1 \leftarrow 1)}\right) = \frac{W_A e^{-\eta}}{W_A e^{-\eta} + W_B},$$

where \mathbb{P} denotes probability with respect to the algorithm's internal randomisation.

6.A.3 Unit Rule Holds for EH on Undirected Spanning Trees

Expanded Hedge on trees can be implemented using the Matrix-Tree Theorem by Kirchoff (undirected) and Tutte (directed). This was pioneered in [99] for log-loss. It can be easily adapted to dot loss. Sampling undirected spanning trees can be done using [23]. This method does not easily generalise to directed spanning trees. Computing the usages is fine for both, and this implies that computing the expected loss is fine for both as well. For directed spanning trees, we can first compute the usages and then decompose as for CH below.

The following theorem neatly characterises the log-partition function.

6.A.1. THEOREM (Kirchhoff's Matrix-Tree Theorem). *Let G be an undirected graph, and let w assign weights to the edges of G . Let \mathcal{S} be the set of*

spanning trees of G , and let L be the graph Laplacian of G , i.e. $L_{i,j} = w(i,j)$ and $L_{i,i} = -\sum_k w(i,k)$. Then

$$\sum_{T \in \mathcal{S}} \prod_{e \in T} w(e) = \det(L_{[1,1]}),$$

where $L_{[1,1]}$ is the first minor of L , i.e. L excluding its first row and column.

This theorem allows us to prove the unit rule for EH on undirected trees.

6.A.2. THEOREM. For all edges e, f

$$Z \cdot Z_{-e \wedge \neg f} \leq Z_{-e} \cdot Z_{\neg f}$$

Simple case (e and f have a common vertex) The following theorem is essential.

6.A.3. THEOREM. For any numbers a, b, c , vectors v, w and symmetric matrix R

$$\det \begin{pmatrix} a & b & v^T \\ b & c & w^T \\ v & w & R \end{pmatrix} \det(R) \leq \det \begin{pmatrix} a & v^T \\ v & R \end{pmatrix} \det \begin{pmatrix} c & w^T \\ w & R \end{pmatrix}$$

Proof. First use the fact that $\det \begin{pmatrix} A & B \\ C & D \end{pmatrix} = \det(D) \det(A - BD^{-1}C)$ repeatedly. We then need to show

$$\begin{aligned} \left(a - (b \ v^T) \begin{pmatrix} c & w^T \\ w & R \end{pmatrix}^{-1} \begin{pmatrix} b \\ v^T \end{pmatrix} \right) (c - w^T R w) \det(R) \det(R) \leq \\ (a - v^T R^{-1} v) \det(R) (c - w^T R w) \det(R) \end{aligned}$$

Now divide out $\det(R)^2$, which is positive, and use

$$\begin{pmatrix} A & B \\ C & D \end{pmatrix}^{-1} = \begin{pmatrix} (A - BD^{-1}C)^{-1} & -(A - BD^{-1}C)^{-1}BD^{-1} \\ -D^{-1}C(A - BD^{-1}C)^{-1} & D^{-1} + D^{-1}B(A - BD^{-1}C)^{-1}CD^{-1} \end{pmatrix}$$

to obtain

$$(b \ v^T) \begin{pmatrix} c & w^T \\ w & R \end{pmatrix}^{-1} \begin{pmatrix} b \\ v^T \end{pmatrix} = \frac{(b - v^T R^{-1} w)^2}{c - w^T R^{-1} w} + v^T R^{-1} v.$$

It then remains to show

$$\left(a - \frac{(b - \mathbf{v}^T R^{-1} \mathbf{w})^2}{c - \mathbf{w}^T R^{-1} \mathbf{w}} - \mathbf{v}^T R^{-1} \mathbf{v} \right) (c - \mathbf{w}^T R \mathbf{w}) \leq (a - \mathbf{v}^T R^{-1} \mathbf{v})(c - \mathbf{w}^T R \mathbf{w})$$

which follows from

$$(b - \mathbf{v}^T R^{-1} \mathbf{w})^2 \geq 0. \quad \square$$

Hard case (e and f have no common vertex)

6.A.4. THEOREM. *Let a, b, c, d, e, f be numbers, $\mathbf{u}, \mathbf{v}, \mathbf{w}$ be vectors, and let R be a symmetric matrix. Then*

$$\begin{aligned} & \left| \begin{array}{cccc} a & b & c & \mathbf{u}^T \\ b & d & e & \mathbf{v}^T \\ c & e & f & \mathbf{w}^T \\ \mathbf{u} & \mathbf{v} & \mathbf{w} & R \end{array} \right| \left(\left| \begin{array}{cc} d & \mathbf{v}^T \\ \mathbf{v} & R \end{array} \right| + \left| \begin{array}{cc} e & \mathbf{v}^T \\ \mathbf{w} & R \end{array} \right| + \left| \begin{array}{cc} e & \mathbf{w}^T \\ \mathbf{v} & R \end{array} \right| + \left| \begin{array}{cc} f & \mathbf{w}^T \\ \mathbf{w} & R \end{array} \right| \right) \leq \\ & \left| \begin{array}{ccc} d & e & \mathbf{v}^T \\ e & f & \mathbf{w}^T \\ \mathbf{v} & \mathbf{w} & R \end{array} \right| \left(\left| \begin{array}{ccc} a & c & \mathbf{u}^T \\ c & f & \mathbf{w}^T \\ \mathbf{u} & \mathbf{w} & R \end{array} \right| + \left| \begin{array}{ccc} a & b & \mathbf{u}^T \\ c & e & \mathbf{w}^T \\ \mathbf{u} & \mathbf{w} & R \end{array} \right| + \left| \begin{array}{ccc} a & c & \mathbf{u}^T \\ b & e & \mathbf{v}^T \\ \mathbf{u} & \mathbf{w} & R \end{array} \right| + \left| \begin{array}{ccc} a & b & \mathbf{u}^T \\ b & d & \mathbf{v}^T \\ \mathbf{u} & \mathbf{v} & R \end{array} \right| \right) \end{aligned}$$

Proof. Currently left to the reader :-). We need a good reduction from this case to the simple case. \square

6.A.4 Unit Rule Holds for FPL on Experts

The unit rule for FPL holds for all perturbations. Fix prior loss ℓ , and let ρ denote a random permutation vector. Then by monotonicity of probability distributions (the right set is bigger)

$$\begin{aligned} \mathbb{P} \left(\bigcap_{k \in [n]} \left\{ \ell_j + \frac{\rho}{\eta} \leq \ell_k + \frac{\rho}{\eta} \right\} \right) &= \mathbb{P}(j = \mathbf{E}^1) \leq \\ \mathbb{P}(j = \mathbf{E}^2 | \ell^1 = \delta_i) &= \mathbb{P} \left(\bigcap_{k \in [n]} \left\{ \ell_j + \frac{\rho}{\eta} \leq (\ell + \delta_i)_k + \frac{\rho}{\eta} \right\} \right). \end{aligned}$$

By the unit rule, we may replace (6.3) by (6.10), eliminating a factor d from under the square root, yielding — up to constants — the same bound as EH. In fact, for Gumbel perturbations FPL coincides with EH. Unfortunately, this fact remains outside the scope of the general analysis of Section 6.2.2.

6.A.5 Unit Rule Holds for FPL on Sets

The usage of component i equals the probability that we draw i when we draw k items without replacement, with probabilities proportional to their weight. Let $0 \leq \beta \leq 1$. Define

$$w_S := \sum_{h \in S} w_h, \quad R(i, S, 0) := 0, \quad R(i, j, S, 0) := 0,$$

and recursively

$$\begin{aligned} R(i, S, k) &:= \frac{w_i + \sum_{h \in S} w_h R(i, S - h, k - 1)}{w_i + w_S} \\ R(i, j, S, k) &:= \frac{w_i + w_j R(i, S, k - 1) + \sum_{h \in S} w_h R(i, S - h, k - 1)}{w_i + w_j + w_S} \\ \tilde{R}(i, j, S, k) &:= \frac{w_i + \beta w_j R(i, S, k - 1) + \sum_{h \in S} w_h \tilde{R}(i, S - h, k - 1)}{w_i + \beta w_j + w_S} \end{aligned}$$

6.A.5. THEOREM. For all $0 \leq k \leq n = |S| + 2$.

$$R(i, j, S, k) = \mathbb{P}(i \in \mathbf{S}^1) \leq \mathbb{P}(i \in \mathbf{S}^2 | \ell^1 = \delta_j) = \tilde{R}(i, j, S, k).$$

Proof. By induction on k . Equality holds for $k = 0$. Suppose the theorem holds up till k . We need to show

$$\begin{aligned} \frac{w_i + w_j R(i, S, k) + \sum_{h \in S} w_h R(i, S - h, k)}{w_i + w_j + w_S} &\leq \\ &\frac{w_i + \beta w_j R(i, S, k) + \sum_{h \in S} w_h \tilde{R}(i, S - h, k)}{w_i + \beta w_j + w_S} \end{aligned}$$

We apply the induction hypothesis, multiply by both denominators, rearrange and divide by $(1 - \beta)w_j(w_i + w_S)$. It then suffices to show

$$R(i, S, k + 1) = \frac{w_i + \sum_{h \in S} w_h R(i, S - h, k)}{w_i + w_S} \geq R(i, S, k)$$

which is obvious. □

6.A.6 Unit Rule Fails for FPL on Permutations

The perturbed loss of a permutation is the loss of that perturbation plus the sum of *two* independent perturbations. Initially both permutations have loss zero, so that either permutation is the perturbed leader with probability one half. If component $(1 \leftarrow 1)$ suffers loss, then obviously the usage of $(2 \leftarrow 2)$ goes down.

6.B Dual Problems for Δ -projection

In this appendix, we compute the Lagrange dual problems to the relative entropy projections on (truncated) permutations, paths and spanning trees. The dual problems involve less variables and less constraints, and thus lead to more efficient implementation of the projection.

6.B.1 Matching Polytope

We want to find

$$\underset{w \text{ s.t. (6.5)}}{\operatorname{argmin}} \Delta(w \|\widehat{w}).$$

Introduce Lagrange multipliers λ_i and μ_j for each row and column constraint. Form the Lagrangian

$$F(w, \lambda, \mu) = \Delta(w \|\widehat{w}) + \sum_{i \in [k]} \lambda_i \left(\sum_{j \in [n]} w_{i \leftarrow j} - 1 \right) + \sum_{j \in [n]} \mu_j \left(\sum_{i \in [k]} w_{i \leftarrow j} - 1 \right)$$

Equating the derivative of F w.r.t. w to zero yields $w^{i \leftarrow j} = \widehat{w}_{i \leftarrow j} e^{-\lambda_i - \mu_j}$. So the dual function is

$$F(\lambda, \mu) := \inf_{w \in \mathbb{R}^d} F(w, \lambda, \mu) = \sum_{i \in [k]} \sum_{j \in [n]} (1 - e^{-\lambda_i - \mu_j}) \widehat{w}_{i \leftarrow j} - \sum_{i \in [k]} \lambda_i - \sum_{j \in [n]} \mu_j$$

The advantage of the dual problem is that we only have $k + n$ variables, whereas the primal has kn . But since μ correspond to inequality constraints, we have to maximise the dual function F under the constraint $\mu \geq \mathbf{0}$.

6.B.2 Flow Polytope

We want to find

$$\operatorname{argmin}_{w \text{ a flow (6.6)}} \Delta(w \|\widehat{w}).$$

Here we generalise slightly. We allow flow from a node to itself. And we allow flow to enter the source. It is still forbidden for flow to leave the sink, and the source still has unit excess flow. We introduce a Lagrange multiplier λ_i , for each node/constraint $i \neq t$, and form the Lagrangian

$$F(w, \lambda) = \Delta(w \|\widehat{w}) + \lambda_s + \sum_{i \neq t} \lambda_i \left(\sum_{j \neq t} w_{j,i} - \sum_j w_{i,j} \right)$$

Equating the derivative of F w.r.t. w to zero yields $w_{i,j} = \widehat{w}_{i,j} e^{\lambda_i - \lambda_j}$, with the convention that $\lambda_t = 0$. The concave dual function thus equals

$$F(\lambda) := \inf_{w \in \mathbb{R}^d} F(w, \lambda) = \lambda_s + \sum_{i \neq t} \sum_j (1 - e^{\lambda_i - \lambda_j}) \widehat{w}_{i,j}.$$

The advantage of the dual problem is that we now only have $n + 1$ variables and no constraints. The primal problem has order n^2 variables and $n + 1$ equality constraints. By strong duality, the optimal primal variables w^* can be reconstructed from the optimum dual variables λ^* .

6.B.3 Tree Polytope

Setup Fix any $\widehat{w} \in \mathbb{R}_+^d$. We are interested in finding

$$\operatorname{argmin}_{w \text{ s.t. (6.7)}} \Delta(w \|\widehat{w}).$$

We introduce Lagrange multipliers α_{ij}^k , γ_{ij}^k , λ and μ_i^k , and form the Lagrangian

$$\begin{aligned} F(w, \alpha, \gamma, \lambda, \mu) := & \Delta(w \|\widehat{w}) - \sum_{i,j \neq i, k \neq 1} \alpha_{ij}^k f_{ij}^k + \sum_{i,j \neq i, k \neq 1} \gamma_{ij}^k (f_{ij}^k - w_{ij}) \\ & + \lambda \left(\sum_{i,j \neq i} w_{ij} - (n-1) \right) + \sum_{i, k \neq 1} \mu_i^k \left(\sum_{j \neq i} (f_{ij}^k - f_{ji}^k) + \theta_i^k \right). \end{aligned}$$

Lagrange dual function The partial derivatives are

$$\frac{\partial F}{\partial w_{ij}} = \log \frac{w_{ij}}{\widehat{w}_{ij}} + \lambda - \sum_{k \neq 1} \gamma_{ij}^k, \quad \frac{\partial F}{\partial f_{ij}^k} = \mu_i^k - \mu_j^k + \gamma_{ij}^k - \alpha_{ij}^k.$$

By setting the partial derivatives to zero, we obtain the Lagrange dual

$$\begin{aligned} F^* &:= \inf_{w_{ij}, f_{ij}^k} F = \sum_{i,j \neq i} w_{ij}^0 - e^{-\lambda} \sum_{i,j \neq i} \widehat{w}_{ij} e^{\sum_{k \neq 1} \gamma_{ij}^k} - \lambda(n-1) + \sum_{i,k \neq 1} \mu_i^k \theta_i^k \\ &= \sum_{i,j \neq i} \widehat{w}_{ij} - e^{-\lambda} \sum_{i,j \neq i} \widehat{w}_{ij} e^{\sum_{k \neq 1} \gamma_{ij}^k} - \lambda(n-1) + \sum_{k \neq 1} (\mu_k^k - \mu_1^k) \end{aligned}$$

Lagrange dual problem We now maximise the Lagrange dual F^* over the dual variables α_{ij}^k , γ_{ij}^k , μ_i^k and λ subject to the constraints

$$\alpha_{ij}^k \geq 0, \quad \gamma_{ij}^k \geq 0, \quad \mu_i^k - \mu_j^k + \gamma_{ij}^k - \alpha_{ij}^k = 0.$$

Since α_{ij}^k do not appear in the dual, these are equivalent to

$$\gamma_{ij}^k \geq 0, \quad \gamma_{ij}^k \geq \mu_j^k - \mu_i^k.$$

Eliminating λ Note that λ is unconstrained. Its derivative is

$$\frac{\partial F^*}{\partial \lambda} = e^{-\lambda} \sum_{i,j \neq i} \widehat{w}_{ij} e^{\sum_{k \neq 1} \gamma_{ij}^k} - (n-1).$$

Setting the derivative to zero, we obtain

$$\begin{aligned} F^\circ &:= \sup_{\lambda} F^* = \\ &\sum_{i,j \neq i} \widehat{w}_{ij} - (n-1) - (n-1) \ln \frac{\sum_{i,j \neq i} \widehat{w}_{ij} e^{\sum_{k \neq 1} \gamma_{ij}^k}}{n-1} + \sum_{k \neq 1} (\mu_k^k - \mu_1^k) \end{aligned}$$

Eliminating γ_{ij}^k Since F° is decreasing in γ_{ij}^k , they each have to be set to their lower bound $\max\{0, \mu_j^k - \mu_i^k\}$. We get

$$\begin{aligned} F^+ &:= \sup_{\gamma_{ij}^k} F^\circ = \sum_{i,j \neq i} \widehat{w}_{ij} - (n-1) - \\ &(n-1) \ln \frac{\sum_{i,j \neq i} \widehat{w}_{ij} e^{\sum_{k \neq 1} \max\{0, \mu_j^k - \mu_i^k\}}}{n-1} + \sum_{k \neq 1} (\mu_k^k - \mu_1^k). \end{aligned}$$

Recovering primal variables Say that we have μ_i^k . We now want to extract the primal variables w_{ij} and f_{ij}^k from these. First we solve for all other dual variables in terms of μ_i^k :

$$\begin{aligned}\gamma_{ij}^k &= \max\{0, \mu_j^k - \mu_i^k\} \\ \alpha_{ij}^k &= \mu_i^k - \mu_j^k + \gamma_{ij}^k = \max\{0, \mu_i^k - \mu_j^k\} = \gamma_{ji}^k \\ \lambda &= \ln \frac{\sum_{i,j \neq i} \widehat{w}_{ij} e^{\sum_{k \neq 1} \gamma_{ij}^k}}{n-1} = \ln \frac{\sum_{i,j \neq i} \widehat{w}_{ij} e^{\sum_{k \neq 1} \max\{0, \mu_j^k - \mu_i^k\}}}{n-1}\end{aligned}$$

We then solve for the primal variables using the KKT conditions

$$\begin{aligned}w_{ij} &= \widehat{w}_{ij} e^{-\lambda + \sum_{k \neq 1} \gamma_{ij}^k} = \frac{(n-1) \widehat{w}_{ij} e^{\sum_{k \neq 1} \max\{0, \mu_j^k - \mu_i^k\}}}{\sum_{i,j \neq i} \widehat{w}_{ij} e^{\sum_{k \neq 1} \max\{0, \mu_j^k - \mu_i^k\}}} \\ f_{ij}^k &= \begin{cases} w_{ij} & \mu_j^k > \mu_i^k, \text{ i.e. } \gamma_{ij}^k > 0 \\ 0 & \mu_j^k < \mu_i^k, \text{ i.e. } \alpha_{ij}^k > 0 \end{cases}\end{aligned}$$