# Bandit Online Optimization Over the Permutahedron

Nir Ailon      Kohei Hatano      Eiji Takimoto

September 13, 2018

### Abstract

The permutahedron is the convex polytope with vertex set consisting of the vectors $(\pi(1), \ldots, \pi(n))$ for all permutations (bijections) $\pi$ over $\{1, \ldots, n\}$. We study a bandit game in which, at each step $t$, an adversary chooses a hidden weight weight vector $s_t$, a player chooses a vertex $\pi_t$ of the permutahedron and suffers an observed instantaneous loss of $\sum_{i=1}^{n} \pi_t(i) s_t(i)$.

We study the problem in two regimes. In the first regime, $s_t$ is a point in the polytope dual to the permutahedron. Algorithm CombBand of Cesa-Bianchi et al (2009) guarantees a regret of $O(n\sqrt{T \log n})$ after $T$ steps. Unfortunately, CombBand requires at each step an $n$-by-$n$ matrix permanent computation, a #P-hard problem. Approximating the permanent is possible in the impractical running time of $O(n^{10})$, with an additional heavy inverse-polynomial dependence on the sought accuracy. We provide an algorithm of slightly worse regret $O(n^{3/2}\sqrt{T})$ but with more realistic time complexity $O(n^3)$ per step. The technical contribution is a bound on the variance of the Plackett-Luce noisy sorting process's 'pseudo loss', obtained by establishing positive semi-definiteness of a family of 3-by-3 matrices of rational functions in exponents of 3 parameters.

In the second regime, $s_t$ is in the hypercube. For this case we present and analyze an algorithm based on Bubeck et al.'s (2012) OSMD approach with a novel projection and decomposition technique for the permutahedron. The algorithm is efficient and achieves a regret of $O(n\sqrt{T})$, but for a more restricted space of possible loss vectors.

## 1 Introduction

Consider a game in which, at each step, a player plays a permutation of some ground set $V = \{1, \ldots, n\}$, and then suffers (and observes) a loss. We model the loss as a sum over the items of some latent quality of the item, weighted by its position in the permutation. The game is repeated, and the items' quality can adversarially change over time. The game models many scenarios in which the player is an online system (say, a search/recommendation engine) presenting a ranked list of items (results/products) to a stream of users. A user's experience

1

is positive if she perceives the quality of the top items on the list as higher than those at the bottom. The goal of the system is to create a total positive experience for its users.

There is a myriad of methods for modelling *ranking* loss functions in the literature, especially (but not exclusively) for information retrieval. Our choice allows us to study the problem in the framework of online combinatorial optimization in the *bandit* setting, and to obtain highly nontrivial results improving on state of the art in either run time or regret bounds. More formally, we study online linear optimization over the the *n-permutahedron* action set, defined as the convex closure of all vectors in $\mathbb{R}^n$ consisting of $n$ distinct coordinates taking values in $[n] := \{1, \ldots, n\}$ (permutations). At each step $t = 1, \ldots, T$, the player outputs an action $\pi_t$ and suffers a loss $\pi_t' s_t = \sum_{i=1}^n \pi_t(i) s_t(i)$, where $s_t \in \mathbb{R}^n$ is the vector of "item qualities" chosen by some adversary who knows the player's strategy but doesn't control their random coins. The performance of the player is the difference between their total loss and that of the optimal static player, who plays the best (in hindsight) single permutation $\pi^*$ throughout. This difference is known as *regret*. Note that, given $s_1, \ldots, s_T$, $\pi^*$ can be computed by sorting the coordinates of $\sum_{t=1}^T s_t$ in decreasing order. This is aligned with our practical requirement that items with higher quality should be placed first, and those with lower quality should be last.

## 2    Results, Techniques and Contribution

Our first of two results, stated as Theorem 1, is for the setting in which at each step the loss is uniformly bounded (by 1 for simplicity) in absolute value for all possible permutations. Equivalently, the vectors $s_t$ belong to the polytope that is dual to the permutahedron. Our algorithm, BanditRank, plays permutations from a distribution known as the Plackett-Luce model (see [12]) which is widely used in statistics and econometrics (see eg [3]). It uses an inverse covariance matrix of the distribution in order to obtain an unbiased loss vector estimator, which is a standard technique [6]. The main technical difficulty (Lemma 2) is in bounding second moment properties of Plackett-Luce, by establishing positive semidefiniteness of a certain family of 3 by 3 matrices. The lemma is interesting in its own right as a tool for studying distributions over permutations. The expected regret of our algorithm is $O(n^{3/2}\sqrt{T})$ for $T$ steps, with running time of $O(n^3)$ per time step. This result should be compared to CombBand of [6], where a framework for playing bandit games over combinatorially structured sets was developed. Their techniques extend that of [7]. In each step, it draws a permutation from a distribution that assigns to each permutation $\pi$ a probability of $e^{\eta \sum_{\tau=1}^t \pi' \tilde{s}_\tau}$, where $\tilde{s}_t$ is a *pseudo-loss* vector at time $t$, an unbiased estimator of the loss vector $s_t$. Their algorithm guarantees a regret of $O(n\sqrt{T \log n})$, which is better than ours by a factor of $\Theta(\sqrt{n/\log n})$. However, its computational requirements are much worse. In order to draw permutations, they need to compute nonnegative $n$ by $n$ matrix permanents. Unfortunately, nonnegative permanent computation is #P-hard, as shown by [14]. On the other hand, a

groundbreaking result of [11] presents a polynomial time approximation scheme for permanent, which runs in time $O(n^{10})$ for fixed accuracy. To make things worse, the dependence in the accuracy is inverse polynomial, implying that, even if we could perform arbitrarily accurate floating point operations, the total running time would be *super linear* in $T$, because a regret dependence of $\sqrt{T}$ over $T$ steps requires accuracy inverse polynomial in $T$. (Our algorithm does not suffer from this problem.) From a practical point of view, the runtime dependence of CombBand in both $n$ and $T$ is infeasible for even modest cases. For example, our algorithm can handle online ranking of $n = 100$ items in an order of few millions of operations per game iteration. In contrast, approximating the permanent of a 100-by-100 positive matrix is utterly impractical.

We note that independently of our work, Hazan et al. [9] have improved the state-of-the-art general purpose algorithm for linear bandit optimization, implying an algorithm with regret $O(n\sqrt{T})$ for our problem, but with worse running time $\tilde{O}(n^4)$.[1]

In our second result in Section 5 we further restrict $s_t$ to have $\ell_1$ norm of $1/n$. (Note that this restriction is contained in $|\pi_t's_t| \le 1$ by Hölder). We present and analyze an algorithm OSMDRank based on the bandit algorithm OSMD of [5] with projection and decomposition techniques over the permutahedron ([15, 13]). The projection is defined in terms of the binary relative entropy divergence. The restriction allows us to obtain an expected regret bound of $O(n\sqrt{T})$ (a $\sqrt{\log n}$ improvement over CombBand). The running time is $O(n^2 + n\tau(n))$, where $\tau(n)$ is the time complexity for some numerical procedure, which is $O(n^2)$ in a fixed precision machine.

We note previous work on playing the permutahedron online optimization game in the *full information case*, namely, when $s_t$ is known for each $t$. As far as we know, Helmbold et al. [10] were the first to study a more general version of this problem, where the action set is the vertex set of the Birkhoff-von-Neumann polytope (doubly-stochastic matrices). Suehiro et al. [13] studied the problem by casting it as a submodularly constrained optimization problem, giving near optimal regret bounds, and more recently Ailon [1] both provided optimal regret bounds with improved running time and established tight regret lower bounds.

## 3 Definitions and Problem Statement

Let $V$ be a ground set of $n$ items. For simplicity, we identify $V$ with $[n] := \{1, \ldots n\}$. Let $S_n$ denote the set of $n!$ permutations over $V$, namely bijections over $[n]$. By convention, we think of $\pi(v)$ for $v \in V$ as the *position* of $v \in V$ in the ranking, where we think of lower numbered positions as *more favorable*. For distinct $u, v \in V$, we say that $u \prec_\pi v$ if $\pi(u) < \pi(v)$ (in words: *u beats v*). We use $[u, v]_\pi$ as shorthand for the indicator function of the predicate $u \prec_\pi v$.

---

[1] The running time is a product of $\tilde{O}(n^3)$ number of Markov chain steps required for drawing a random point from a convex set under a log-concave distribution, and $O(n \log n)$ time to test whether a point lies in the permutahedron. By $\tilde{O}$ we hide poly-logarithmic factors.

The convex closure of $S_n$ is known as the permutahedron polytope. It will be more convenient for us to consider a translated version of the permutahedron, centered around the origin. More precisely, for $\pi \in S_n$ we let $\hat{\pi}$ denote

$$\hat{\pi} := (\pi(1) - (n+1)/2,\, \pi(2) - (n+1)/2, \ldots, \pi(n) - (n+1)/2) \ .$$

It will be convenient to define a symmetrized version of the permutation set $\hat{S}_n := \{\hat{\pi} : \pi \in S_n\}$. The symmetrized $n$-permutahedron, denoted $\hat{P}_n$ is the convex closure of $\hat{S}_n$. Symmetrization allows us to work with a polytope that is centered around the origin. Generalization our result to standard (un-symmetrized) permutations is a simple technicality that will be explained below. The notation $u \prec_{\hat{\pi}} v$ and $[u,v]_{\hat{\pi}}$ is defined as for $\pi \in S_n$ in an obvious manner.

At each step $t = 1, \ldots, T$, an adversary chooses and hides a nonnegative vector $s_t \in \mathbb{R}^n \equiv \mathbb{R}^V$, which assigns an elementwise quality measure $s_t(v)$ for any $v \in V$. The player-algorithm chooses a permutation $\hat{\pi}_t \in \hat{S}_n$, possibly random, and suffers an instantaneous loss

$$\ell_t := \hat{\pi}_t' s_t = \sum_{v \in V} \hat{\pi}_t(v) s_t(v) \ . \tag{3.1}$$

The total loss $L_t$ is defined as $\sum_{t=1}^{T} \ell_t$. We will work with the notion of *regret*, defined as the difference $L_t - L_t^*$, where $L_T^* = \min_{\hat{\pi} \in \hat{S}_n} \sum_{t=1}^{T} \hat{\pi}' s_t$. We let $\hat{\pi}^*$ denote any minimizer achieving $L_T^*$ in the RHS.

For any $\hat{\pi} \in \hat{S}_n$ and $s \in \mathbb{R}^n$, the dot-product $\hat{\pi}' s$ can be decomposed over pairs: $\hat{\pi}' s = \frac{1}{2} \sum_{u \neq v} [u,v]_\pi (s(v) - s(u))$. This makes the symmetrized permutahedron easier to work with. Nevertheless, our results also apply to the non-symmetrized permutahedron as well, as we shall see below.

Throughout, the notation $\sum_{u \neq v}$ means summation over distinct, ordered pairs of elements $u, v \in V$, and $\sum_{u < v}$ means summation over distinct, unordered pairs.[2] The uniform distribution over $\hat{S}_n$ will be denoted $\mathcal{U}_n$.

The smallest eigenvalue of a PSD matrix $A$ is denoted $\lambda_{\min}(A)$. The norm $\|\cdot\|_2$ will denote spectral norm (Euclidean norm for a vector). To avoid notation such as $C, C', C'', C_1$ for universal constants, the expression $C$ will denote a "general positive constant" that may change its value as necessary. For example, we may write $C = 3C + 5$.

# 4 Algorithm BanditRank and its Guarantee

For this section, we will assume that the instantaneous losses are uniformly bounded by 1, in absolute value: For all $t$ and $\hat{\pi} \in \hat{S}_n$, $|\hat{\pi}' s_t| \leq 1$. Equivalently, using geometric language, the loss vectors belong to a polytope which is *dual* to the permutahedron.

Now consider Algorithm 1. It maintains, at each time step $t$, a weight vector $w_t \in \mathbb{R}^n$. At each time step, it draws a random permutation $\hat{\pi}_t$ from a mixture

---

[2] We will only use expressions of the form $\sum_{u<v} f(u,v)$ for symmetric functions satisfying $f(u,v) = f(v,u)$.

$\mathcal{D}_t$ of the uniform distribution over $\hat{S}_n$ and a distribution $\mathcal{PL}_n(w)$ which we define shortly. The distribution mixture is determined by a parameter $\gamma$. The algorithm then plays the permutation $\hat{\pi}_t$ and thereby suffers the instantaneous loss defined in (3.1). The weights are consequently updated by adding an unbiased estimator $\tilde{s}_t$ of $s_t$ (computed using the pseudo-inverse covariance matrix corresponding to $\mathcal{D}_t$), multiplied by another parameter $\eta > 0$.

**The Plackett-Luce Random Sorting Procedure:** The distribution $\mathcal{PL}_n(w)$ over $\hat{S}_n$, parametrized by $w \in \mathbb{R}^n$, is defined by the following procedure. To choose the first (most preferred) item, the procedure draws a random item, assigning probability proportional to $e^{w(u)}$ for each $u \in V$. It then removes this item from the pool of available items, and iteratively continues to choose the second item, then third and so on. As claimed in the introduction, this random permutation model is well studied in statistics. An important well known property of the distribution is that it can be equivalently defined as a *Random Utility Model (RUM)* [12, 16]: To draw a permutation, add a random iid noise variable following the Gumbel distribution to each weight, and then sort the items of $V$ in decreasing value of noisy-weights.[3] The RUM characterization implies, in particular, that for any two *disjoint* pairs of element $(u, v)$ and $(u', v')$, the events $u \prec_\pi v$ and $u' \prec_\pi v'$ are statistically independent if $\pi$ is drawn from $\mathcal{PL}_n(w)$, for any $w$. This fact will be used later.
We are finally ready to state our main result, bounding the expected regret of the algorithm.

**Theorem 1.** *If algorithm* BanditRank *(Algorithm 1) is executed with parameters $\gamma = O(n^{3/2}/\sqrt{T})$ and $\eta = O(\gamma/n)$, then the expected regret (with respect to the game defined by the symmetrized permutahedron) is at most $O(n^{3/2}\sqrt{T})$. The running time of each iteration is $O(n^3)$. Additionally, there exists an algorithm with the same expected regret bound and running time with respect to the standard permutahedron (assuming the vectors $s_t$ uniformly satisfy $|\pi' s_t| \le 1, \forall \pi \in S_n$.)*

The proof uses a standard technique used e.g. in Cesa-Bianchi et al.'s Comb-Band [6], which is itself an adaptation of Auer et al.'s Exp3 [2] from the finite case to the structured combinatorial case. The distribution from which the actions $\hat{\pi}_t$ are drawn in the algorithm differ from the distribution used in Comb-Band, and give rise to the technical difficulty of variance estimation, resolved in Lemma 2.

*Proof.* Let $\mathcal{T}_n$ denote the set of *tournaments* over $[n]$. More precisely, an element $A \in \mathcal{T}_n$ is a subset of $[n] \times [n]$ with either $(u, v) \in A$ or $(v, u) \in A$ (but not both) for all $u < v$. We extend our previous notation so that $u \prec_A v$ is equivalent to the predicate $(u, v) \in A$.

---

[3] The Gumbel distribution, also known as doubly-exponential, has a cdf of $e^{-e^{-x}}$.

**Algorithm 1** Algorithm BanditRank($n, \eta, \gamma, T$) (assuming $|\hat{\pi}'s_t| \leq 1$ for all $t$ and $\hat{\pi} \in \hat{S}_n$)

---

1: given: ground set size $n$, positive parameters $\eta, \gamma$ ($\gamma \leq 1$), time horizon $T$
2: set $w_0(u) = 0$ for all $u \in V = [n]$
3: **for** $t = 1..T$ **do**
4:     let distribution $\mathcal{D}_t$ over $\hat{S}_n$ denote a mixture of $\mathcal{U}_n$ (with probability $\gamma$) and $\mathcal{PL}_n(w_{t-1})$ (with probability $1 - \gamma$)
5:     draw and output $\hat{\pi}_t \sim \mathcal{D}_t$
6:     observe and suffer loss $\ell_t$ ($= \hat{\pi}_t's_t$)
7:     $\tilde{s}_t = \ell_t P_t^+ \hat{\pi}_t$ where $P_t = \mathbb{E}_{\hat{\sigma} \sim \mathcal{D}_t}[\hat{\sigma}\hat{\sigma}']$
8:     set $w_t = w_{t-1} + \eta\tilde{s}$
9: **end for**

---

For any pair $\hat{\pi} \in \hat{S}_n$ and $w \in \mathbb{R}^n$, $p(\hat{\pi}|w)$ denotes the probability assigned to $\hat{\pi} \in \hat{S}_n$ by $\mathcal{PL}_n(w)$. Slightly abusing notation, we define the following shorthand:

$$p(u \prec v|w) := \sum_{\hat{\pi}:u \prec_{\hat{\pi}} v} p(\pi|w) = \frac{e^{w(u)}}{e^{w(u)} + e^{w(v)}}$$

$$p(u \prec v \prec z|w) := \sum_{\hat{\pi}:u \prec_{\hat{\pi}} v \prec_{\hat{\pi}} z} p(\hat{\pi}|w) = \frac{e^{w(u)+w(v)}}{(e^{w(u)} + e^{w(v)} + e^{w(z)})(e^{w(v)} + e^{w(z)})} .$$

The last two right hand sides are easily derived from the definition of the distribution $\mathcal{PL}_n(w)$, see also e.g. [12]. We also define the following abbreviations:

$$p(u \prec {}^v_z|w) := p(u \prec v \prec z|w) + p(u \prec z \prec v|w) = \frac{e^{w(u)}}{e^{w(u)} + e^{w(v)} + e^{w(z)}} \tag{4.1}$$

$$p({}^u_v \prec z|w) := p(u \prec v \prec z|w) + p(v \prec u \prec z|w)$$
$$= \frac{e^{w(u)+w(v)}}{e^{w(u)} + e^{w(v)} + e^{w(z)}} \left( \frac{1}{e^{w(v)} + e^{w(z)}} + \frac{1}{e^{w(u)} + e^{w(z)}} \right) \tag{4.2}$$

We will also need to define a distribution over the set of tournaments $\mathcal{T}_n$. The distribution, $\mathcal{BTL}_n(w)$ is parametrized by a weight vector $w \in \mathbb{R}^n$. Drawing $A \sim \mathcal{BTL}_n(w)$ is done by independently setting, for all $u < v$ in $V$,

$$(u, v) \in A \text{ with probability } p(u \prec v|w) = \frac{e^{w(u)}}{e^{w(u)} + e^{w(v)}}$$
$$(v, u) \in A \text{ with probability } p(v \prec u|w) = \frac{e^{w(v)}}{e^{w(u)} + e^{w(v)}} .$$

(Note that the distribution is equivalently defined as the product distribution, over all $u < v$ in $V$, of the Bradley-Terry-Luce pairwise preference model, hence

the name $\mathcal{BTL}_n$. We refer to [12] for definition and history of the Bradley-Terry-Luce model.)

For $A \in \mathcal{T}_n$, we denote by $\tilde{p}(A|w)$ the probability $\prod_{u \prec_A v} p(u \prec v|w)$ of drawing $A$ from $\mathcal{BTL}_n(w)$. The proof of the theorem proceeds roughly as the main result upper bounding the expected regret of CombBand in [6]. The following technical lemma is required in anticipation of a major hurdle (inequality (4.5). We believe the inequality is interesting in its own right as a probabilistic statement on permutation and tournament distributions.

**Lemma 2.** *Let $s, w \in \mathbb{R}^n$. Let $\hat{\pi} \sim \mathcal{PL}_n(w)$ and $A \sim \mathcal{BTL}_n(w)$ be drawn independently. Define $X_1 = \sum_{u,v:\ u \prec_{\hat{\pi}} v}(s(v) - s(u)) = \hat{\pi}'s$, $X_2 = \sum_{u,v:\ u \prec_A v}(s(v) - s(u))$. Then $\mathbb{E}[X_2^2] \leq \mathbb{E}[X_1^2]$.*

(Note that clearly, $\mathbb{E}[X_2] = \mathbb{E}[X_1]$, so the lemma in fact upper bounds the variance of $X_2$ by that of $X_1$.) The proof of the lemma is deferred to Section 4.1.

Continuing the proof of Theorem 1, we let $q(\pi|w)$ denote the probability of drawing $\pi$ from the mixture of the uniform distribution (with probability $\gamma$) and $\mathcal{PL}_n(w)$ (with probability $(1 - \gamma)$). Similarly to above, $q(u \prec v|w)$ denotes $\sum_{\hat{\pi}:u \prec_{\hat{\pi}} v} q(\hat{\pi}|w)$. By these definitions,

$$q(\hat{\pi}|w) = (1 - \gamma)p(\hat{\pi}|w) + \frac{\gamma}{n!} \qquad q(u \prec v|w) = (1 - \gamma)p(u \prec v|w) + \frac{\gamma}{2} \ . \quad (4.3)$$

The analysis proceeds by defining a potential function: $W_t(u, v) := e^{\frac{1}{2}\eta(w_t(u) - w_t(v))} + e^{\frac{1}{2}\eta(w_t(v) - w_t(u))}$. The quanatity of interest will be $\mathbb{E}\left[\sum_{u<v}\sum_t \log \frac{W_t(u,v)}{W_{t-1}(u,v)}\right]$, where the expectation is taken over all random coins used by the algorithm throughout $T$ steps. This quantity will be bounded from above and from below, giving rise to a bound on the expected total loss, expressed using the optimal static loss. On the one hand,

$$\sum_{u<v} \log \frac{W_t(u,v)}{W_{t-1}(u,v)} = \sum_{u<v} \log \left( \frac{e^{\frac{1}{2}(w_t(u) - w_t(v))}}{W_{t-1}(u,v)} + \frac{e^{\frac{1}{2}(w_t(v) - w_t(u))}}{W_{t-1}(u,v)} \right)$$

$$= \sum_{u<v} \log \left( \frac{e^{\frac{1}{2}(w_{t-1}(u) - w_{t-1}(v))}e^{\frac{1}{2}\eta(\tilde{s}_t(u) - \tilde{s}_t(v))}}{W_{t-1}(u,v)} + \frac{e^{\frac{1}{2}(w_t(v) - w_t(u))}e^{\frac{1}{2}\eta(\tilde{s}_t(v) - \tilde{s}_t(u))}}{W_{t-1}(u,v)} \right)$$

$$= \sum_{u<v} \log \left( p(u \prec v|w_{t-1})e^{\frac{1}{2}\eta(\tilde{s}_t(u) - \tilde{s}_t(v))} + p(v \prec u|w_{t-1})e^{\frac{1}{2}\eta(\tilde{s}_t(v) - \tilde{s}_t(u))} \right)$$

$$= \log \left( \sum_{A \in \mathcal{T}_n} \tilde{p}(A|w_{t-1})e^{\frac{1}{2}\eta \sum_{u \prec_A v}(\tilde{s}_t(u) - \tilde{s}_t(v))} \right) \ .$$

We will now assume that $\eta$ is small enough so that for all $A \in \mathcal{T}_n$ and for all $t$,

$$\eta \left| \sum_{(u,v) \in A} (\tilde{s}_t(u) - \tilde{s}_t(v)) \right| \leq 1 \ . \tag{4.4}$$

7

(This will be shortly enforced.) Using $e^x \leq 1 + x + x^2 \ \forall x \in [-1/2, 1/2]$,

$$\sum_{u,v} \log \frac{W_t(u,v)}{W_{t-1}(u,v)} \leq \log \left[ \sum_{A \in \mathcal{T}_n} \tilde{p}(A|w_{t-1}) \left( 1 + \frac{\eta}{2} \sum_{u \prec_A v} (\tilde{s}_t(u) - \tilde{s}_t(v)) \right. \right.$$

$$\left. \left. + \frac{\eta^2}{4} \left( \sum_{u \prec_A v} (\tilde{s}_t(u) - \tilde{s}_t(v)) \right)^2 \right) \right]$$

$$= \log \left[ 1 + \frac{\eta}{2} \mathbb{E}_{A \sim \mathcal{BTL}_n(w_{t-1})} \left[ \sum_{u \prec_A v} (\tilde{s}_t(u) - \tilde{s}_t(v)) + \frac{\eta^2}{4} \left( \sum_{u \prec_A v} (\tilde{s}_t(u) - \tilde{s}_t(v)) \right)^2 \right] \right]$$

$$\leq \log \left[ 1 + \frac{\eta}{2} \mathbb{E}_{\hat{\pi} \sim \mathcal{PL}_n(w_{t-1})} \left[ \sum_{u \prec_{\hat{\pi}} v} (\tilde{s}_t(u) - \tilde{s}_t(v)) + \frac{\eta^2}{4} \left( \sum_{u \prec_{\hat{\pi}} v} (\tilde{s}_t(u) - \tilde{s}_t(v)) \right)^2 \right] \right] .$$

$$(4.5)$$

where we used Lemma 2 in the last inequality (together with the fact that the marginal probability of the event "$u \prec_Y v$" is identical for both $Y \sim \mathcal{PL}_n(w_{t-1})$ and $Y \sim \mathcal{BTL}_n(w_{t-1})$). Henceforth, for any $\hat{\pi} \in \hat{S}$, we let $\tilde{\ell}_t(\hat{\pi}) := \hat{\pi}' \tilde{s}_t = \sum_{u \prec_{\hat{\pi}} v} (\tilde{s}(v) - \tilde{s}(u))$. Using 4.3 and the fact that $\log(1+x) \leq x$ for all $x$, we get

$$\sum_{u < v} \log \frac{W_t(u,v)}{W_{t-1}(u,v)}$$

$$\leq \frac{\eta}{2} \sum_{u \neq v} \frac{q(u \prec v|w_{t-1}) - \frac{\gamma}{2}}{1 - \gamma} (\tilde{s}_t(u) - \tilde{s}_t(v)) + \frac{\eta^2}{4} \sum_{\hat{\pi} \in \hat{S}_n} \frac{q(\pi|w_{t-1}) - \frac{\gamma}{n!}}{1 - \gamma} \tilde{\ell}_t(\hat{\pi})^2$$

$$\leq \frac{-\eta}{2(1-\gamma)} \sum_{\hat{\pi} \in \hat{S}_n} q_t(\hat{\pi}|w_{t-1}) \tilde{\ell}_t(\hat{\pi}) + \frac{\eta^2}{4(1-\gamma)} \sum_{\hat{\pi} \in \hat{S}_n} q_t(\hat{\pi}|w_{t-1}) \tilde{\ell}_t(\hat{\pi})^2 .$$

We now note that (1) $\sum_{\hat{\pi} \in \hat{S}} q_t(\hat{\pi}|w_{t-1}) \tilde{\ell}_t = \ell_t$ (following the properties of matrix pseudo-inverse in Line 7 in Algorithm 1), and (2) $\sum_{\hat{\pi} \in \hat{S}_n} q_t(\hat{\pi}|w_{t-1}) \tilde{\ell}_t(\pi)^2] \leq n$ (see top of page 31 together with Lemma 15 in [6]). Applying these inequalities, and then taking expectations over the algorithm's randomness and summing for $t = 1, \ldots, T$, we get

$$\sum_{t=1}^{T} \mathbb{E} \left[ \sum_{u,v} \log \frac{W_t(u,v)}{W_{t-1}(u,v)} \right] \leq -\frac{\eta}{2(1-\gamma)} \mathbb{E}[L_T] + \frac{\eta^2}{8(1-\gamma)} nT .$$

8

On the other hand,

$$\sum_{t=1}^{T}\mathbb{E}\left[\sum_{u,v}\log\frac{W_t(u,v)}{W_{t-1}(u,v)}\right]$$

$$\geq \sum_{u,v}\mathbb{E}\left[\log\left([u,v]_{\pi^*}e^{\frac{1}{2}(w_T(u)-w_T(v))}+[v,u]_{\pi^*}e^{\frac{1}{2}(w_T(u)-w_T(v))}\right)\right)\right]-\sum_{u,v}\log 2$$

$$=\frac{1}{2}\sum_{u,v}\left(\mathbb{E}\left[[u,v]_{\pi^*}(w_T(u)-w_T(v))+[v,u]_{\pi^*}(w_T(u)-w_T(v))\right]\right)-\binom{n}{2}\log 2$$

$$=\frac{\eta}{2}\sum_{u,v}\left(\mathbb{E}\left[[u,v]_{\pi^*}\sum_t(\tilde{s}_t(u)-\tilde{s}_t(v))+[v,u]_{\pi^*}\sum_t(\tilde{s}_t(u)-\tilde{s}_t(v))\right]\right)-\binom{n}{2}\log 2$$

$$=\frac{\eta}{2}\sum_{u,v}\left([u,v]_{\pi^*}\sum_t(s_t(u)-s_t(v))+[v,u]_{\pi^*}\sum_t(s_t(u)-s_t(v))\right)-\binom{n}{2}\log 2$$

$$=-\frac{\eta}{2}L_T^*-\binom{n}{2}\log 2\;,$$

where $L_T^*$ is the total loss of a player who chooses the best permutatation $\hat{\pi}^* \in \hat{S}_n$ in hindsight. Combining, we obtain $\frac{\eta}{2(1-\gamma)}\mathbb{E}[L_t] \leq \frac{\eta}{2}L_T^*+\frac{n^2}{2}\log 2+\frac{\eta^2}{4(1-\gamma)}nT$. Multiplying both sides by $2(1-\gamma)/\eta$ yields

$$\mathbb{E}[L_T] \leq L_T^* + \gamma|L_T^*| + \frac{n^2\log 2}{\eta} + \frac{\eta}{2}nT\;. \tag{4.6}$$

We shall now work to impose (4.4).

$$\max_t\max_{A\in\mathcal{T}(V)}\left|\sum_{(u,v)\in A}(\tilde{s}_t(u)-\tilde{s}_t(v))\right| \leq \max_t\sqrt{\sum_{v\in V}\tilde{s}_t(v)^2}\sqrt{\sum_{i=-(n-1)/2}^{(n-1)/2}i^2} \leq C\max_t\|\tilde{s}_t\|_2 n^{3/2}\;,$$

where the left inequality is Cauchy-Schwartz. We now note that $\|\tilde{s}_t\|_2 \leq |\ell_t|\|P_t^+\|_2\|\hat{\pi}_t\|_2$. Clearly $\|\hat{\pi}\|_2$ is bounded above by $Cn^{3/2}$. Also $\|P_t^+\|_2$ equals $1/\lambda_{\min}(P_t)$. By Weyl's inequality $\lambda_{\min}(P_t) \geq \gamma\lambda_{\min}(\mathbb{E}_{\hat{\tau}\sim\mathcal{U}_n}[\hat{\tau}\hat{\tau}'])$. It is an exercise to check that $\lambda_{\min}(\mathbb{E}_{\hat{\tau}\sim\mathcal{U}_n}[\hat{\tau}\hat{\tau}']) \geq Cn^2$. We conclude (also recalling that $|\ell_t| \leq 1$) that $\max_t\|\tilde{s}_t\|_2 \leq C/(n^{1/2}\gamma)$. Combining, we shall satisfy (4.7) by imposing $\eta \leq \gamma/(Cn)$. Plugging in (4.6), we get

$$\mathbb{E}[L_T(\mathrm{Alg})] \leq L_T^* + \gamma|L_T^*| + \frac{Cn^3}{\gamma} + C\gamma T\;. \tag{4.7}$$

Choosing $\gamma = \sqrt{\frac{Cn^3}{T}}$ gives $\mathbb{E}[L_T(\mathrm{Alg})] \leq L_T^* + \frac{Cn^{3/2}}{\sqrt{T}}|L_T^*| + n^{3/2}\sqrt{T}$.

This concludes the required result for the symmetrized case, because $|L_T^*| \leq T$. For the standard permutahedron, we notice that for any $\pi \in S_n$ and its symmetrized counterpart $\hat{\pi} \in \hat{S}_n$, and any vector $s \in \mathbb{R}^n$, $\pi's - \hat{\pi}'s =$

9

$\frac{n-1}{2} \sum_{v \in V} s(v) =: f(s)$. Equivalently, we can write $\pi's = (\hat{\pi}', 1)(s; f(s))$, where $(\cdot, a)$ appends the scalar $a$ to the right of a row vector and $(\cdot; a)$ appends to the bottom of a column vector. Algorithm 1 can be easily adjusted to work with action set $\hat{S}_n \times \{1\}$. For the proof, we keep the same potential function. The technical part of the proof is lower bounding the smallest eigenvalue of the expectation of $\hat{\tau}\hat{\tau}'$, where $\hat{\tau}$ is now drawn from the uniform distribution on $\hat{S}_n \times \{1\}$. We omit these simple details for lack of space. $\qquad\square\qquad\qquad\square$

## 4.1 Proof of Lemma 2

The expression $\mathbb{E}[X_1^2]$ can be written as

$$
\begin{aligned}
\mathbb{E}[X_1^2] = &\sum_{u \neq v} p(u \prec v | w)((s(v) - s(u))^2 \\
&+ \sum_{|\{u,v,u',v'\}|=4} p(u \prec v \wedge u' \prec v'|w)(s(v) - s(u))(s(v') - s(u')) \\
&+ \sum_{\substack{u \neq v, u' \neq v' \\ |\{u,v,u',v'\}|=3}} p(u \prec v \wedge u' \prec v'|w)(s(v) - s(u))(s(v') - s(u')) , \quad (4.8)
\end{aligned}
$$

where $p(u \prec v \wedge u' \prec v'|w)$ is the probability that both $u \prec_{\hat{\pi}} v$ and $u' \prec_{\hat{\pi}} v'$ with $\hat{\pi} \sim \mathcal{PL}_n(w)$. Similarly,

$$
\begin{aligned}
\mathbb{E}[X_2^2] = &\sum_{u \neq v} p(u, v | w)((s(v) - s(u))^2 \\
&+ \sum_{|\{u,v,u',v'\}|=4} p(u \prec v|w)p(u' \prec v'|w)(s(v) - s(u))(s(v') - s(u')) \\
&+ \sum_{\substack{u \neq v, u' \neq v' \\ |\{u,v,u',v'\}|=3}} p(u \prec v|w)p(u' \prec v'|w)(s(v) - s(u))(s(v') - s(u')) .
\end{aligned}
$$

$$(4.9)$$

Since Plackett-Luce is a random utility model (see [12]), it is clear that whenever a pair of pairs $u \neq v, u' \neq v'$ satisfies $|\{u, v, u', v'\}| = 4$, $p(u \prec v \wedge u' \prec v'|w) = p(u \prec v|w)p(u' \prec v'|w)$. Hence, it suffices to prove that the third summand in the RHS of (4.9) is upper bounded by the third summand in the RHS of (4.8). But now notice the following identity:

$$
\sum_{\substack{u \neq v, u' \neq v' \\ |\{u,v,u',v'\}|=3}} \equiv \sum_{\substack{\Delta \subseteq V \\ |\Delta|=3}} \sum_{\substack{u \neq v, u' \neq v' \\ u,v,u',v' \in \Delta \\ |\{u,v,u',v'\}|=3}} .
$$

This last sum rearrangement implies that it suffices to prove that for any $\Delta$ of

cardinality 3,

$$F_2(\Delta) := \sum_{\substack{u \neq v, u' \neq v' \\ u,v,u',v' \in \Delta \\ |\{u,v,u',v'\}|=3}} p(u,v|w)p(u',v'|w)\,(s(v)-s(u))(s(v')-s(u'))$$

$$\leq \sum_{\substack{u \neq v, u' \neq v' \\ u,v,u',v' \in \Delta \\ ||\{u,v,u',v'\}|=3}} p(u, v \wedge u', v'|w)\,(s(v)-s(u))(s(v')-s(u')) =: F_1(\Delta) \ .$$

If we now denote $\Delta = \{a, b, c\}$, then both $F_1(\Delta)$ and $F_2(\Delta)$ are quadratic forms in $s(a), s(b), s(c)$ (for fixed $w$). It hence suffices to prove that $H(\Delta) := F_1(\Delta) - F_2(\Delta)$ is a positive semi-definite form in $s(\Delta) := (s(a), s(b), s(c))'$. We now write

$$H(\Delta) = s(\Delta)' \begin{pmatrix} H_{aa} & \frac{1}{2}H_{ab} & \frac{1}{2}H_{ac} \\ \frac{1}{2}H_{ab} & H_{bb} & \frac{1}{2}H_{bc} \\ \frac{1}{2}H_{ac} & \frac{1}{2}H_{bc} & H_{cc} \end{pmatrix} s(\Delta) \ .$$

The matrix is singular, because clearly $H(\Delta) = F_1(\Delta) = F_2(\Delta) = 0$ whenever $s(a) = s(b) = s(c)$. To prove positive semi-definiteness, by Sylvester's criterion it hence suffices to show that the diagonal element $H_{aa} \geq 0$ and that the principal 2-by-2 minor determinant $H_{aa}H_{bb} - \frac{1}{4}H_{ab}^2 \geq 0$. Using the definitions, together with the properties of $\mathcal{PL}_n(w)$, a technical (but quite tedious) algebraic derivation (see Appendix A for details) gives

$$H_{aa} = \frac{4e^{s(a)+s(b)+s(c)}}{(e^{s(a)} + e^{s(b)})(e^{s(a)} + e^{s(c)})(e^{s(a)} + e^{s(b)} + e^{s(c)})} \ . \tag{4.10}$$

Similarly, by symmetry, $H_{bb} = \frac{4e^{s(a)+s(b)+s(c)}}{(e^{s(b)}+e^{s(a)})(e^{s(b)}+e^{s(c)})(e^{s(a)}+e^{s(b)}+e^{s(c)})}$. From a similar (yet more tedious) technical algebraic calculation which we omit, one gets: (see Appendix A for details):

$$H_{ab} = \frac{-8e^{s(a)+s(b)+2s(c)}}{(e^{s(a)} + e^{s(b)})(e^{s(a)} + e^{s(c)})(e^{s(b)} + e^{s(c)})(e^{s(a)} + e^{s(b)} + e^{s(c)})} \ . \tag{4.11}$$

One now verifies, using (4.10)-(4.11), the identity

$$H_{aa}H_{bb} - \frac{1}{4}H_{ab}^2 = \frac{16e^{2s(a)+2s(b)+2s(c)}}{(e^{s(a)} + e^{s(b)})^2(e^{s(a)} + e^{s(c)})(e^{s(b)} + e^{s(c)})(e^{s(a)} + e^{s(b)} + e^{s(c)})^2} \ .$$

It remains to notice, trivially, that $H_{aa} \geq 0$ and $H_{aa}H_{bb} - \frac{1}{4}H_{ab}^2 \geq 0$ for all possible values of $s(a), s(b), s(c)$. The proof of the lemma is concluded.

# 5 Bandit Algorithm based on Projection and Decomposition

In this section, we propose another bandit algorithm OSMDRank, described in Algorithm 2. We will be working under the more restricted assumption that

**Algorithm 2** Algorithm OSMDRank($n, \eta, \gamma, T$) (assuming $\|s_t\|_1 \leq 1$ and $\hat{\pi}_t \in \hat{Q}_n$ for all $t$ )

---

1: given: ground set size $n$, positive parameters $\eta, \gamma$ ($\gamma \leq 1$), time horizon $T$
2: let $x_1 = 0 \in \hat{Q}_n$. (Note that $x_1 = \arg\min_{a \in \hat{Q}_n} F(a)$)
3: **for** $t = 1, \ldots, T$ **do**
4:     let $\tilde{x}_t = (1 - \gamma)x_t$ (Note that $\tilde{a}_t \in \hat{Q}_n$ since the origin 0 and $x_t$ are in $\hat{Q}_n$ and $\tilde{x}_t$ is a convex combination of them).
5:     output $\pi_t = \text{Decomposition}(\tilde{x}_t)$ (i.e., choose $\pi_t$ so that $\mathbb{E}[\pi_t] = \tilde{x}_t$) and suffer loss $\ell_t$ ($= \pi_t's_t$)
6:     let distribution $\mathcal{D}_t$ over $[-1, 1]^n$ denote a mixture of the uniform distribution over the canonical basis with random sign (with probability $\gamma$) and a Radmacher distribution over $\{-1, 1\}^n$ with parameter $(1 + x_{t,i})/2$ for each $i = 1, \ldots, n$ (with probability $1 - \gamma$)
7:     estimate the loss vector $\tilde{s}_t = \ell_t P_t^+ \pi_t$, where $P_t = \mathbb{E}_{\sigma \sim D_t}[\sigma\sigma']$
8:     let $x_{t+\frac{1}{2}} = \nabla F^*(F(x_t) - \eta\tilde{s}_t)$
9:     let $x_{t+1} = \text{Projection}(x_{t+\frac{1}{2}})$ (that is, $x_{t+1} = \min_{x \in \hat{Q}_n} D_F(x, x_{t+\frac{1}{2}})$)
10: **end for**

---

$\sup \|s_t\|_1 \leq 1$ and $\sup \|\hat{\pi}_t\|_\infty \leq 1$. This in particular implies that $|\hat{\pi}_t's_t| \leq 1$, as before. But now we shall achieve a better expected regret of $O(n\sqrt{T})$.

We prefer, for reasons clarified shortly, to require that the actions $\hat{\pi}_t$ are vertices of the rescaling $\hat{Q}_n := \frac{2}{n-1}\hat{P}_n \in [-1, 1]^n$ of the symmetrized permutahedron. That is, $\sup \|\hat{\pi}_t\|_\infty \leq 1$ (and $\sup \|s_t\|_1 \leq 1$). This will allow us to work with the following standard regularizer $F : [-1, 1]^n \to \mathbb{R}^+$: $F(x) = \frac{1}{2}\sum_{i=1}^n ((1+x)\ln(1+x) + (1-x)\ln(1-x))$. The regularizer $F(x)$ is the key to the OSMD (Online Stochastic Mirror Descent) algorithm of Bubeck et al. [5], on which our algorithm is based. OSMD is a bandit algorithm over the hypercube domain $[-1, 1]^n$ and a variant of Follow the Regularized Leader (FTRL, e.g., [8]) for linear loss functions. To apply this algorithm, we need a new projection and decomposition technique for the polytope $\hat{Q}_n$, as well as a slightly modified perturbation step in line 4 of Algorithm 2. Our algorithm OSMDRank has the following two procedures:

1. **Projection:** Given a point $x_t \in [-1, 1]^n$, return $\arg\min_{y_t \in \hat{Q}_n} \Delta_F(y_t, x_t)$, where $\Delta_F$ is the Bregman divergence defined wr.t. $F$, i.e., $\Delta_F(y, x) = F(y) - F(x) - \nabla F(x)'(y - x)$ (also known as *binary relative entropy*).[4]

2. **Decomposition:** Given $y_t \in \hat{Q}_n$ from the the projection step, output a random vertex $\hat{\pi}_t$ of $\hat{Q}_n$ such that $\mathbb{E}[\hat{\pi}_t] = y_t$.

The decomposition can be done using the technique of [15], which runs in $O(n \log n)$ time. (To be precise, the method there was defined for the standard

---

[4]Note that the binary relative entropy is different from the relative entropy, where the relative entropy is defined as $Rel(p, q) = \sum_{i=1}^n p_i \ln \frac{p_i}{q_i}$ for probability distributions $p$ and $q$ over $[n]$.

permutahedron; The adjustments for the symmetrized version are trivial.) For notational purposes, we define $f := \nabla F$, and notice that $f(x)_i = \frac{1}{2} \ln \frac{1+x_i}{1-x_i}$, and its inverse function $f^{-1}$ is given by $f^{-1}(y)_i = \frac{e^{y_i}-1}{e^{y_i}+1}$. Our projection procedure is presented in Algorithm 3.

**Lemma 3.** *(i) Given $q \in [-1,1]^n$, Algorithm 3 outputs the projection of $q$ onto $\hat{Q}_n$, with respect to the regularizer $F$. (ii) The time complexity of the algorithm is $O(n\tau(n) + n^2)$, where $\tau(n)$ is the time complexity to perform step 4.*

*skecth.* Our projection algorithm is an extension of that in [13] and our proof follows a similar argument in [13]. For simplicity, we assume that elements in $q$ are sorted in descending order, i.e., $q_1 \geq q_2 \geq \cdots \geq q_n$. This can be achieved in time $O(n \log n)$ by sorting $q$. Then, it can be shown that projection preserves the order in $q$ by using Lemma 1 in [13]. That is, the projection $p$ of $q$ satisfies $p_1 \geq p_2 \geq \cdots \geq p_n$. So, if the conditions $\frac{2}{n-1} \sum_{j=1}^{i} p_j \leq \sum_{j=1}^{i} (\frac{n+1}{2} - j)$, for $i = 1, \ldots, n-1$, are satisfied, then other inequality constraints are satisfied as well since for any $S \subset [n]$ such that $|S| = i$, $\sum_{j \in S} p_j \leq \sum_{j=1}^{i} p_j$. Therefore, relevant constraints for projection onto $\hat{Q}_n$ are only linearly many.

By following a similar argument in [13], we can show that the output $p$ indeed satisfies the KKT optimality conditions for projection, which completes the proof of the first statement. Finally, the algorithm terminates in time $O(n\tau(n) + n^2)$ since the number of iteration is at most $n$ and each iteration takes $O(n + \tau(n))$ time, which completes the second statement of the lemma. $\square$ $\square$

Note that with respect to other regularizers (e.g. relative entropy or Euclidean norm squared), a different projection scheme is possible in time $O(n^2)$ (see [15, 13] for the details). It is an open question whether an $O(n^2)$ algorithm can be devised with respect to the binary relative entropy we need here. In our case, we need to solve a numerical optimization problem by, say, binary search. Note that the time $\tau(n)$ is reasonably small: In fact, we can perform the binary search over the domain $[-1,1]$ for each dimension $i$. Therefore, if the precision is a fixed constant, the binary search ends in time $O(n)$ for each dimension. In that case, $\tau(n)$ is $O(n^2)$. We are ready to present our main result for this section.

**Theorem 4.** *For $\eta = O(n\sqrt{1/T})$ and $\gamma = O(\sqrt{1/T})$, Algorithm OSMDRank has expected regret $O(n\sqrt{T})$ and running time $O(n^2 + n\tau(n))$ per step, where $\tau(n)$ is the time for a numerical optimization step depending on $n$. Additionally, there exists an algorithm with the same expected regret bound and running time with respect to the standard permutahedron (assuming $\|s_t\|_1 \leq 1/n$).*

*sketch.* The algorithm OSMDRank is a modification of OSMD for the hypercube $[-1,1]^n$ obtained by adding (1) a projection step and (2) a decomposition step. Standard techniques show that adding the projection step does not increase the expected regret bound (see, e.g., chapters 5 and 7 on OMD and OSMD of Bubeck's lecture notes [4]). The key facts are: (i) A variant of Theorem 2 of [5] (regret bound of OSMD) holds for OSMD with Projection, (ii) $E[\pi_t] = (1-\gamma)x_t$,

---

**Algorithm 3** Projection onto $\hat{Q}_n$

---

1: given $(q_1, \ldots, q_n) \in [-1, 1]^n$ satisfying $q_1 \geq q_2 \geq \cdots \geq q_n$. *(This assumption holds by renaming the indices, and reverting to their original names at the end).*

2: set $i_0 = 0$

3: **for** $k = 1, \ldots, n$ **do**

4:     for each $i = i_{k-1} + 1, \ldots, n$, set $\delta_i^k = \min_{\delta \in \mathbb{R}} \delta$ subject to:
$$\sum_{j=i_{k-1}+1}^{i} f^{-1}(f(q_j) - \delta) \leq \frac{2}{n-1} \sum_{j=i_{k-1}+1}^{i} \left( \frac{n+1}{2} - j \right).$$

5:     $i_k = \arg\max_{i:i_{k-1} < i \leq n} \delta_i^k$. In case of multiple minimizers, choose largest as $i_k$.

6:     set $p_j = f^{-1}(f(q_j) - \delta_{i_k}^k)$ for $j = i_{k-1} + 1, \ldots, i_k$

7:     **if** $i_k = n$, **then** break

8: **end for**

9: **return** $(p_1, \ldots, p_n)'$

---

and (iii) The estimated loss is the same one used in OSMD for the hypercube $[-1, 1]^n$. Once these three conditions are satisfied, we can prove a regret bound of OSMDRank by following the proof of Theorem 5 in Bubeck et al. [5]. In addition, the running time of OSMD per trial is $O(n)$ [5]. Combining Lemma 3 for the projection and the analysis of the decomposition from [15], the proof of the first statement is concluded. The statement related to the standard permutahedron holds based on the affine transformation between the standard permutahedron and $\hat{Q}_n$.    □                                     □

# 6   Future Work

The main open question is whether there is an algorithm of expected regret $O(n\sqrt{T})$ and time $O(n^3)$ in the setting of Section 4. Another interesting line of research is to study other ranking polytopes. For example, given any strictly monotonically increasing function $f : \mathbb{R} \mapsto \mathbb{R}$ we can consider as an action set $f^n(S_n)$, defined as $f^n(S_n) := \{(f(\pi(1)), f(\pi(2)), \ldots, f(\pi(n))) : \pi \in S_n\}$.

# Acknowledgments

# References

[1] Nir Ailon. Improved Bounds for Online Learning Over the Permutahedron and Other Ranking Polytopes. In *AISTATS*, 2014.

[2] Peter Auer, Nicolò Cesa-Bianchi, Yoav Freund, and Robert E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM J. Comput.*, 32(1):48–77, January 2003.

[3] S Beggs, S Cardell, and J Hausman. Assessing the potential demand for electric cars. *Journal of Econometrics*, 17(1):1 – 19, 1981.

[4] Sébastien Bubeck. Introduction to Online Optimization. http://www.princeton.edu/~bubeck/BubeckLectureNotes.pdf, 2011.

[5] Sébastien Bubeck, Nicolò Cesa-Bianchi, and Sham M. Kakade. Towards Minimax Policies for Online Linear Optimization with Bandit Feedback. In *Proceedings of 25th Annual Conference on Learning Theory (COLT 2012)*, pages 41.1–41.14, 2012.

[6] Nicolò Cesa-Bianchi and Gábor Lugosi. Combinatorial bandits. *J. Comput. Syst. Sci.*, 78(5):1404–1422, 2012.

[7] Varsha Dani, Thomas P. Hayes, and Sham Kakade. The price of bandit information for online optimization. In *NIPS*, 2007.

[8] Elad Hazan. The convex optimization approach to regret minimization. In Suvrit Sra, Sebastian Nowozin, and Stephen J. Wright, editors, *Optimization for Machine Learning*, chapter 10, pages 287–304. MIT Press, 2011.

[9] Elad Hazan, Zohar Shay Karnin, and Raghu Mehka. Volumetric spanners and their applications to machine learning. *CoRR*, abs/1312.6214, 2013.

[10] David P. Helmbold and Manfred K. Warmuth. Learning Permutations with Exponential Weights. *Journal of Machine Learning Research*, 10:1705–1736, 2009.

[11] Mark Jerrum, Alistair Sinclair, and Eric Vigoda. A polynomial-time approximation algorithm for the permanent of a matrix with nonnegative entries. *J. ACM*, 51(4):671–697, 2004.

[12] John I. Marden. *Analyzing and Modeling Rank Data*. Chapman & Hall, 1995.

[13] Daiki Suehiro, Kohei Hatano, Shuji Kijima, Eiji Takimoto, and Kiyohito Nagano. Online Prediction under Submodular Constraints. In *Proceedings of 23th Annual Conference on Algorithmic Learning Theory (ALT 2012)*, pages 260–274, 2012.

[14] Leslie G. Valiant. The complexity of computing the permanent. *Theor. Comput. Sci.*, 8:189–201, 1979.

[15] Shota Yasutake, Kohei Hatano, Shuji Kijima, Eiji Takimoto, and Masayuki Takeda. Online Linear Optimization over Permutations. In *Proceedings of the 22nd International Symposium on Algorithms and Computation (ISAAC 2011)*, pages 534–543, 2011.

[16] J. Yellott. The relationship between Luce's choice axiom, Thurstone's theory of comparative judgment, and the double exponential distribution. *Journal of Mathematical Psychology*, 15:109–144, 1977.

# A    Derivations in proof of Lemma 2

By definition, and then by applying the properties of the distribution $\mathcal{PL}_n(w)$,

$$
\begin{aligned}
H_{aa} = {} & [p(a \prec b \wedge a \prec c|w) + p(b \prec a \wedge c \prec a|w) - p(b \prec a \wedge a \prec c|w) - p(c \prec a \wedge a \prec b|w)] \\
& - [p(a \prec b|w)p(a \prec c|w) + p(b \prec a|w)p(c \prec a|w) - p(a \prec b|w)p(c \prec a|w) \\
& \qquad\qquad\qquad\qquad\qquad - p(a \prec c|w)p(b \prec a|w)] \qquad \text{(A.1)}
\end{aligned}
$$

$$
p(a \prec b \wedge a \prec c|w) = \frac{e^{s(a)}}{e^{s(a)} + e^{s(b)} + e^{s(c)}} \tag{A.2}
$$

$$
p(b \prec a \wedge c \prec a|w) = \frac{e^{s(b)}}{e^{s(a)} + e^{s(b)} + e^{s(c)}} \frac{e^{s(c)}}{e^{s(a)} + e^{s(c)}} + \frac{e^{s(c)}}{e^{s(a)} + e^{s(b)} + e^{s(c)}} \frac{e^{s(b)}}{e^{s(a)} + e^{s(b)}} \tag{A.3}
$$

$$
p(b \prec a \wedge a \prec c|w) = \frac{e^{s(b)}}{e^{s(a)} + e^{s(b)} + e^{s(c)}} \frac{e^{s(a)}}{e^{s(a)} + e^{s(c)}} \tag{A.4}
$$

$$
p(c \prec a \wedge a \prec b|w) = \frac{e^{s(c)}}{e^{s(a)} + e^{s(b)} + e^{s(c)}} \frac{e^{s(a)}}{e^{s(a)} + e^{s(b)}} \tag{A.5}
$$

Plugging (A.2)-(A.5) in (A.1) and simplifying results in (4.10). One now verifies:

$$
\begin{aligned}
H_{ab} = {} & [p(a \prec c \wedge b \prec c|w) + p(c \prec a \wedge c \prec b|w) - 3p(a \prec c \wedge c \prec b|w) - 3p(b \prec c \wedge c \prec a|w)] \\
& - [-p(a \prec b|w)p(a \prec c|w) - p(b \prec a|w)p(c \prec a|w) + p(a \prec b|w)p(c \prec a|w) \\
& \quad + p(a \prec b|w)p(a \prec c|w) + p(a \prec b|w)p(b \prec c|w) + p(b \prec a|w)p(c \prec b|w) \\
& \quad - p(b \prec a|w)p(b \prec c|w) - p(a \prec b|w)p(c \prec b|w) + p(a \prec b|w)p(b \prec c|w) \\
& \quad + p(b \prec a|w)p(c \prec b|w) - p(b \prec a|w)p(b \prec c|w) - p(a \prec b|w)p(c \prec b|w) \\
& \quad - p(a \prec c|w)p(b \prec c|w) - p(c \prec a|w)p(c \prec b|w) + p(a \prec c|w)p(c \prec b|w) \\
& \qquad\qquad\qquad\qquad\qquad\qquad + p(c \prec a|w)p(b \prec c|w)]
\end{aligned}
$$

Again using identities (A.2)-(A.5) and simplifying, gives (4.11)

16