

364 A Deferred proofs

365 A.1 Proof of Corollary 1

366 Consider two potential “bad” events: B_1 , where $\|\hat{\mu}_n - \mu\|_\infty > 1/2$, and B_2 , where $\|\mu - \hat{\mu}_n\|_1 >$
 367 $\frac{2\|\hat{\mu}_n\|_1^{1/2}}{\sqrt{n}} + 6\sqrt{\frac{\ln(4/\delta)}{2n}}$. Our assumption on the sample size n , together with the Dvoretzky-Kiefer-
 368 Wolfowitz inequality [Massart, 1990], implies that $\mathbb{P}(B_1) \leq \delta/2$ and (1) implies that $\mathbb{P}(B_2) \leq \delta/2$.
 369 Thus, with probability at least $1 - \delta$, neither of B_1 or B_2 occurs, and we may invoke Theorem 1,
 370 from which the claim immediately follows.

371 □

372 A.2 Proof of Theorem 2

373 The following non-trivial fact [Lieb and Loss, 2001, Theorem 3.5 and Eq. (5) on p. 83] will be
 374 useful⁶:

$$\|\mu^\downarrow - \nu^\downarrow\|_p \leq \|\mu - \nu\|_p, \quad p \in [1, \infty], \quad \mu, \nu \in \Delta_{\mathbb{N}}. \quad (11)$$

375 A result of Scheffé [1947] (more accurately credited to Riesz, 1928 [Kusolitsch, 2010]) implies
 376 that a sequence $\{\xi_{n \in \mathbb{N}}\} \subset \ell_1(\mathbb{N})$ converging pointwise to some $\xi \in \ell_1(\mathbb{N})$ also converges in ℓ_1 iff
 377 $\|\xi_n\|_1 \rightarrow \|\xi\|_1$. This immediately implies

378 **Lemma 3.** *If $\{\mu_{n \in \mathbb{N}}\} \subset \Delta_{\mathbb{N}}$ converges pointwise to some $\mu \in \Delta_{\mathbb{N}}$, then it also converges in ℓ_1 .*

379 Berend et al. [2017, Lemma 1] showed that $\Delta_{\mathbb{N}}^{\downarrow(1)}[h]$ is compact under ℓ_1 . We begin by extending
 380 this result to general α, p .

381 **Lemma 4.** *For all $\alpha \geq 1, p \in [1, \infty]$, and $h > 0$, the set $\Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$ is compact under ℓ_p .*

382 **Remark.** This is quite false if either the non-increasing or the bounded-entropy condition is omitted.
 383 For a counterexample to the former, consider the sequence $\mu_n \in \Delta_{\mathbb{N}}$ defined by $\mu_n(i) = \mathbf{1}[i = n]$.
 384 For a counterexample to the latter, consider the sequence $\mu_n \in \Delta_{\mathbb{N}}$, where μ_n is uniform on $[n]$.

385 *Proof.* We closely follow the proof strategy of Berend et al. [2017, Lemma 1]. In a metric space,
 386 compactness and sequential compactness are equivalent. Let $\mu_{n \in \mathbb{N}}$ be a sequence in $\Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$. Since
 387 $[0, 1]$ is compact, every $\{\mu_n(i) : n \in \mathbb{N}\}$ has a convergent subsequence, and hence $\mu_{n \in \mathbb{N}}$ has a
 388 pointwise convergent subsequence. There is thus no loss of generality in assuming that $\mu_n \rightarrow \mu$
 389 pointwise. Obviously, μ is non-negative and non-increasing. It remains to show that

390 (a) $\sum_{i \in \mathbb{N}} \mu(i) = 1,$

391 (b) $H^{(\alpha)}(\mu) \leq h,$

392 (c) $\|\mu_n - \mu\|_p \rightarrow 0.$

393 To show (a), assume, for a contradiction, that $\sum_{i \in \mathbb{N}} \mu(i) > 1$. Then there must be an $i_0 \in \mathbb{N}$ such that
 394 $\sum_{i=1}^{i_0} \mu(i) > 1$. But the latter must then hold for all μ_n with n sufficiently large, which contradicts
 395 $\mu_n \in \Delta_{\mathbb{N}}$. Now assume $\varepsilon := 1 - \sum_{i \in \mathbb{N}} \mu(i) > 0$. For any $i_0 \in \mathbb{N}$, we have $\sum_{i=1}^{i_0} \mu_n(i) < 1 - \varepsilon/2$
 396 for all sufficiently large n . Now every $\nu \in \Delta_{\mathbb{N}}^{\downarrow}$ satisfies $\nu(i) \leq \frac{1}{i}(\nu(1) + \nu(2) + \dots + \nu(i)) \leq \frac{1}{i}$.
 397 Hence,

$$\sum_{i=i_0+1}^{\infty} \mu_n(i) |\log \mu_n(i)|^\alpha \geq \sum_{i=i_0+1}^{\infty} \mu_n(i) (\log i_0)^\alpha > \frac{\varepsilon}{2} (\log i_0)^\alpha.$$

⁶The result is stated for functions in $f \in L_2(\mathbb{R}^n)$ and their symmetric-decreasing rearrangements f^* , but the specialization to discrete distributions is straightforward. We convert μ to a function $f : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ via $f(x) = \mu(\lceil x \rceil)$ and ν to $g(x)$ analogously. A direct calculation then shows that $\|\mu - \nu\|_p = \|f - g\|_p$ and $\|\mu^\downarrow - \nu^\downarrow\|_p = \|f^* - g^*\|_p$, to which the result from Lieb and Loss [2001] applies to yield (11).

398 Choosing i_0 sufficiently large makes the latter expression exceed h , violating the assumption $\mu_n \in$
 399 $\Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$. Thus (a) holds.

400 To show (b), assume for a contradiction that $H^{(\alpha)}(\mu) > h$ — and, in particular,
 401 $\sum_{i=1}^{i_0} \mu(i) |\log \mu(i)|^\alpha > h$ for some $i_0 \in \mathbb{N}$. But the latter must hold for all μ_n with n suffi-
 402 ciently large, a contradiction.

403 Finally, to show (c), we invoke Lemma 3: if $\{\mu_{n \in \mathbb{N}}\} \subset \Delta_{\mathbb{N}}$ converges pointwise to some $\mu \in \Delta_{\mathbb{N}}$,
 404 then it also converges in ℓ_1 . Since ℓ_1 dominates every ℓ_p , $p > 1$, this proves (c). \square

405 Next, we examine the continuity of $H(\cdot)$ on $\Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$ under ℓ_p .

406 **Lemma 5.** Fix $h > 0$, $\alpha > 1$, and $p \in [1, \infty]$. If $\{\mu_{n \in \mathbb{N}}\} \subset \Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$ converges in ℓ_p , then its limit
 407 is some $\mu \in \Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$ and furthermore, $H(\mu_n) \rightarrow H(\mu)$. In other words, $H(\cdot)$ is continuous on
 408 $\Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$ under ℓ_p .

409 **Remark.** We note that $H(\cdot)$ is not continuous on $\Delta_{\mathbb{N}}^{\downarrow(1)}[h]$ under ℓ_p , $p \in [1, \infty]$, as evidenced by
 410 the sequence $\mu_n = (1 - \varepsilon_n, \varepsilon_n/n, \dots, \varepsilon/n, 0, 0, \dots)$, with support size $n + 1$. We can choose ε_n so
 411 that $H(\mu_n) = h$, but of course the limiting μ has $H(\mu) = 0$ (see Example 1 in Berend et al. [2017]).

412 *Proof.* It follows from Lemma 4 that the limiting μ belongs to $\Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$. Further, Lemma 3 implies
 413 that $\mu_n \rightarrow \mu$ in ℓ_1 . Invoking the continuity result in Theorem 1 proves the claim. \square

414 *Proof of Theorem 2.* It follows from Lemma 5 that $H(\cdot)$ is continuous on $\Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$ under ℓ_p . Since,
 415 by Lemma 4, $\Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$ is compact under ℓ_p , it follows that $H(\cdot)$ is uniformly continuous on $\Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$:
 416 there is a function F such that

$$|H(\mu) - H(\nu)| \leq F(\|\mu - \nu\|_p, h, \alpha, p), \quad \mu, \nu \in \Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$$

417 and $\varepsilon_n := \|\mu_n - \nu_n\|_p \rightarrow 0 \implies F(\varepsilon_n, h, \alpha, p) \rightarrow 0$. Now, for all $\mu, \nu \in \Delta_{\mathbb{N}}^{\downarrow(\alpha)}[h]$ we have

$$\begin{aligned} |H(\mu) - H(\nu)| &= |H(\mu^\downarrow) - H(\nu^\downarrow)| \\ &\leq F(\|\mu^\downarrow - \nu^\downarrow\|_p, h, \alpha, p). \end{aligned}$$

418 It follows from (11) that $\|\mu_n - \nu_n\|_p \rightarrow 0 \implies \|\mu^\downarrow - \nu^\downarrow\|_p \rightarrow 0$, which concludes the proof. \square

419 B Auxiliary results

420 **Proposition 2.** For $K \geq 2$ and $\alpha \geq 1$,

$$\max \{\log K, (\alpha/e)\}^\alpha \leq \max_{\mu \in \Delta_K} H^{(\alpha)}(\mu) \leq \max \{\alpha, \log K\}^\alpha + (\alpha/e)^\alpha.$$

421 We will need the following useful (and likely known) result.

422 **Lemma 6** (folklore). Suppose that $0 < a < 1$ and $f : [0, 1] \rightarrow \mathbb{R}$ is strictly concave on $[0, a]$ and
 423 strictly convex on $[a, 1]$. Define the function $F : \Delta_K \rightarrow \mathbb{R}$ by

$$F(\mu) = \sum_{i=1}^K f(\mu(i)).$$

424 Then any maximizer μ^* of F is either the uniform distribution or else has exactly 1 “heavy” mass
 425 $v \in [a, 1]$ and $K - 1$ identical “light” masses $(1 - v)/(K - 1)$.

426 *Proof.* A standard “smoothing” argument [Loh, 2013] shows that if two masses $u \leq v$ occur in the
 427 interval $(a, 1)$, there is an $\varepsilon > 0$ such that $f(u - \varepsilon) + f(v + \varepsilon) > f(u) + f(v)$. In other words,
 428 such masses can be pushed apart (keeping their sum fixed) to increase the value of F , until one of
 429 them reaches the boundary of $[a, 1]$. Furthermore, since $0 < a < u < v$ and $u + v \leq 1$, repeated

430 iteration of the “pushing apart” operation will hit the left endpoint (i.e., a) rather than the right one
 431 (i.e., 1). Having exhausted the “pushing apart” process, we are left with one “heavy” mass $v \in [a, 1]$
 432 and $K - 1$ “lighter” ones in $[0, a]$. But concavity implies that F will be maximized by pulling the
 433 lighter masses in (as opposed to pushing them apart), which amounts to replacing each of them by
 434 the average of the $K - 1$ values. \square

435 *Proof of Proposition 2.* Choosing μ to be the uniform distribution yields $H^{(\alpha)}(\mu) = \log^\alpha K$, and
 436 choosing μ such that $v := \mu(1) = e^{-\alpha}$ yields $H^{(\alpha)}(\mu) \geq v \log(1/v)^\alpha = (\alpha/e)^\alpha$. Thus, the lower
 437 bound is proven and it only remains to prove the upper bound.

438 Let μ^* be a maximizer for given α, K . Recall the function $h^{(\alpha)}(z) = z \log^\alpha(1/z)$ and note that it is
 439 strictly concave on $[0, e^{-(\alpha-1)}]$ and strictly convex on $[e^{-(\alpha-1)}, 1]$. Then Lemma 6 shows that μ^*
 440 will either be uniform or else attains at most one value $v \in [e^{-(\alpha-1)}, 1]$ in the convex interval, with
 441 the remaining values equal to $\frac{1-v}{K-1} \in [0, e^{-(\alpha-1)}]$ in the concave interval. Only the latter case is
 442 non-trivial:

$$H^{(\alpha)}(\mu^*) = v \left(\log \frac{1}{v} \right)^\alpha + (1-v) \left(\log \frac{K-1}{1-v} \right)^\alpha$$

443 for some v satisfying

$$0 < \frac{1-v}{K-1} \leq e^{-(\alpha-1)} \leq v < 1. \quad (12)$$

444 Now $v \left(\log \frac{1}{v} \right)^\alpha$ is maximized over $[0, 1]$ by $v = e^{-\alpha}$, which yields the value $(\alpha/e)^\alpha$.

445 To bound the second term, $g(v) := (1-v) \left(\log \frac{K-1}{1-v} \right)^\alpha$, we consider two cases: (i) $K - 1 < e^\alpha$
 446 and (ii) $K - 1 \geq e^\alpha$. In case (i), g is maximized by $v^* = 1 - (K - 1)/e^\alpha$ and

$$g(v^*) = (1-v^*) \left(\log \frac{K-1}{1-v^*} \right)^\alpha \leq \left(\log \frac{K-1}{1-v^*} \right)^\alpha = \alpha^\alpha.$$

447 In case (ii), g is monotonically decreasing in v . The constraint $\frac{1-v}{K-1} \leq e^{-(\alpha-1)}$ from (12) implies $v \geq 1/K$,
 448 so in this case,

$$g(v) \leq \left(\log \frac{K-1}{1-1/K} \right)^\alpha = \log^\alpha K.$$

449 This proves the upper bound.

450 \square